

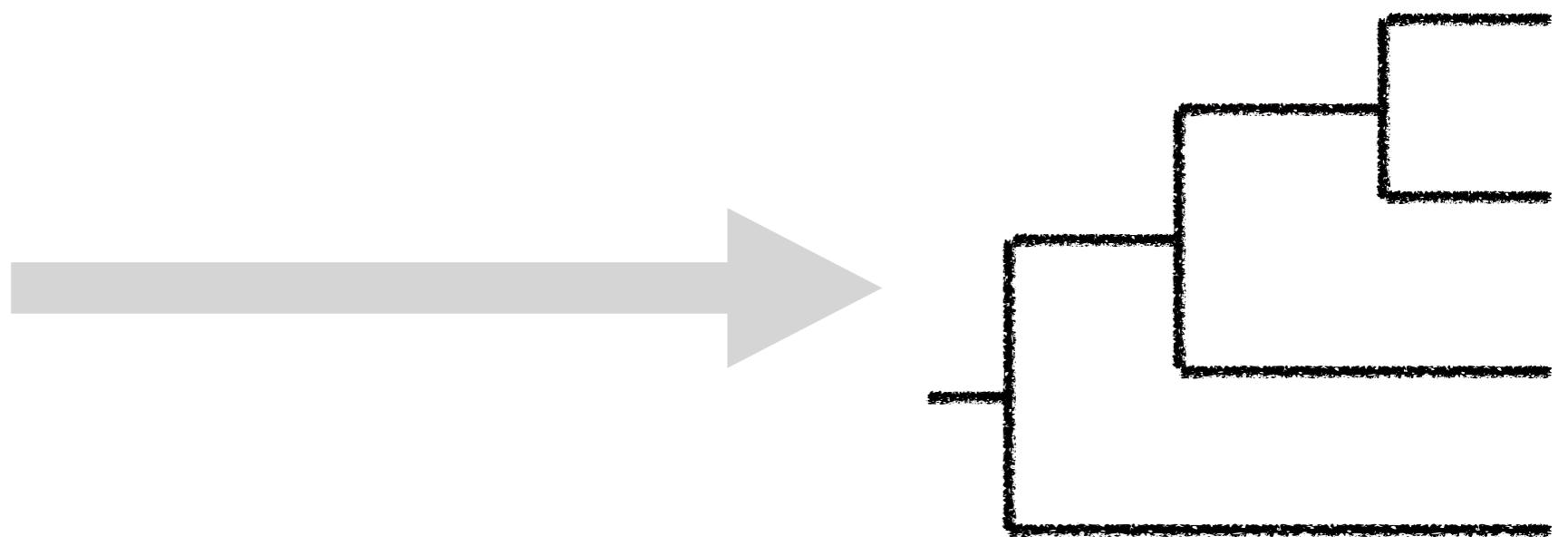
# Lecture 14

Coalescent-based methods  
Botany 563 – Spring 2021

- **Previous class check-up:**
  - We studied Bayesian phylogenetic inference
  - We practiced on MrBayes and/or PhyloBayes
- **Learning Objectives:** At the end of today's session, you will be able to
  - Explain the coalescent model on a species tree
  - Explain the coalescent model on a species network
  - Explain the steps in coalescent-based methods and the comparison with concatenation approaches
- **Pre-class work**
  - Read HAL 3.1 and 3.3

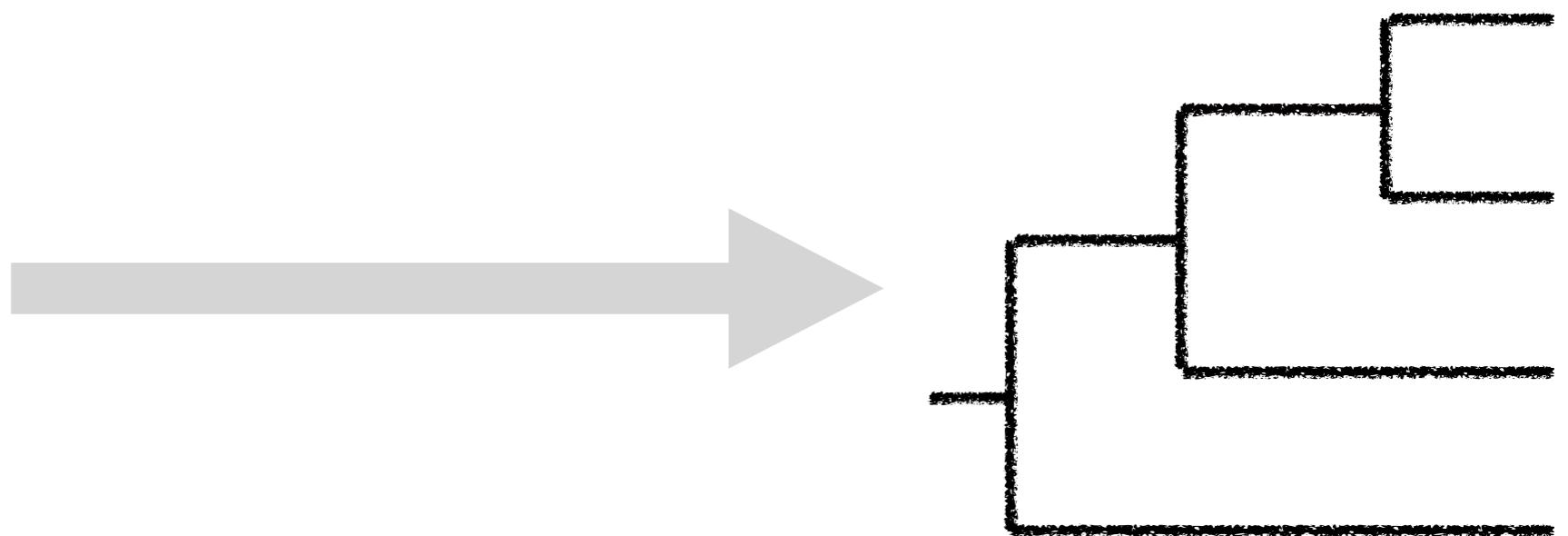
# Phylogenetic inference

AAGTCTAG  
AAGTCTAG  
AACTCTAG  
AATTCTAG

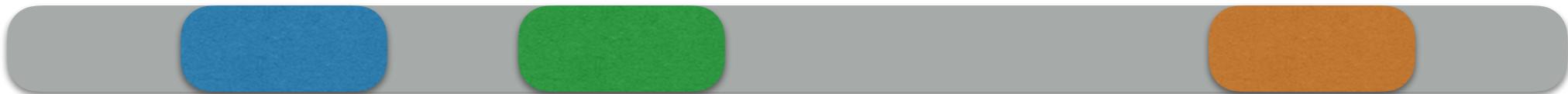


# Phylogenetic inference

AAGTCTAG  
AAGTCTAG  
AACTCTAG  
AATTCTAG



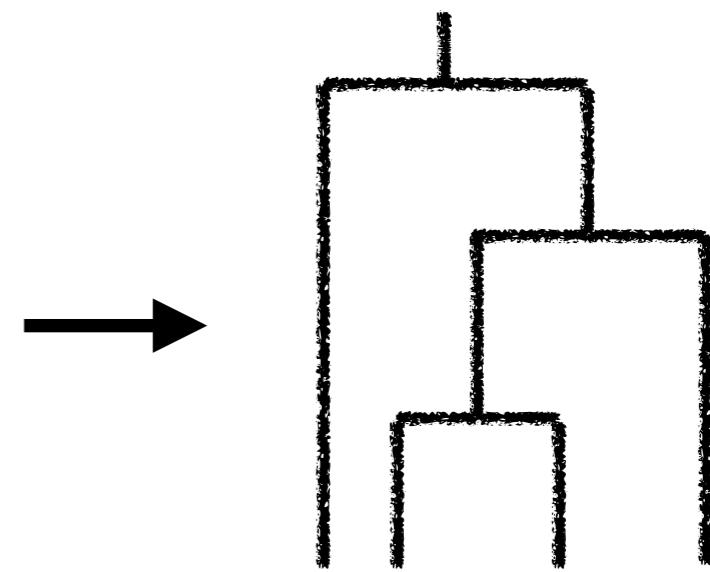
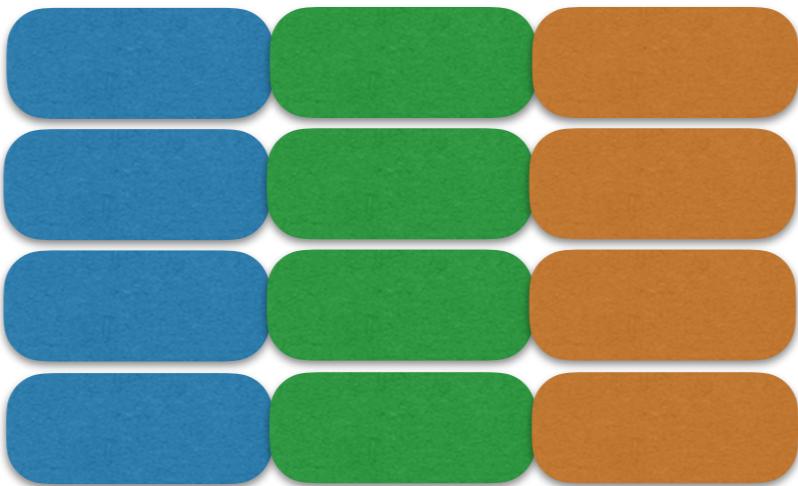
Gene  
Locus  
Region  
Whole genome



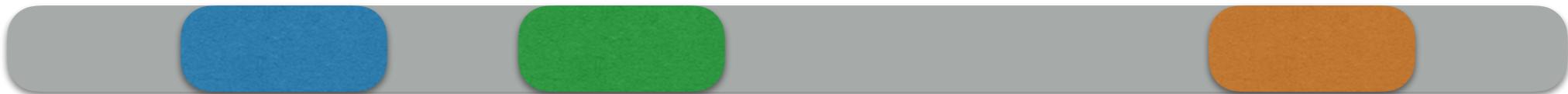


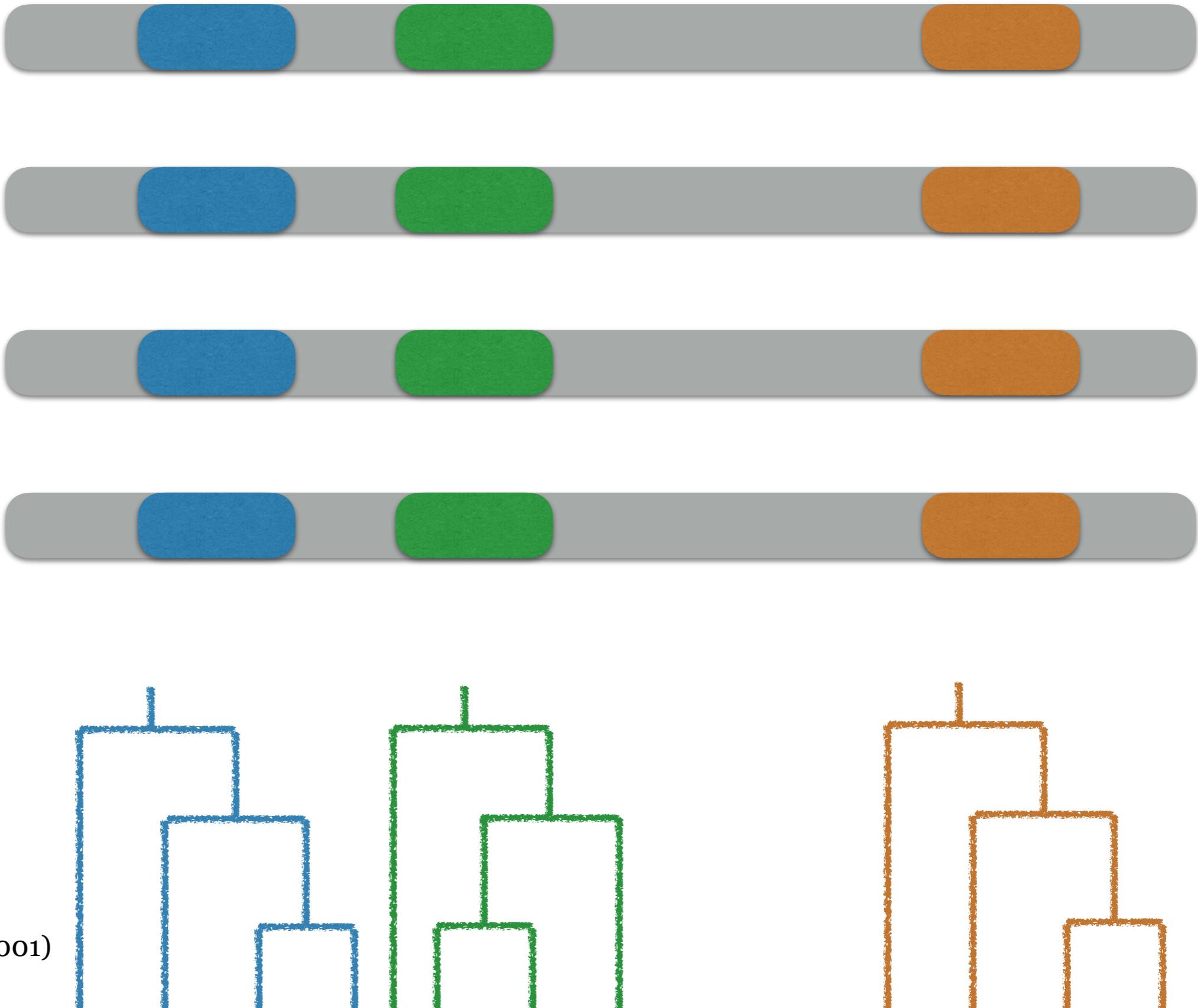
Statistically  
inconsistent

(Kubatko, Degnan, 2007)  
(Roch, Steel, 2015)

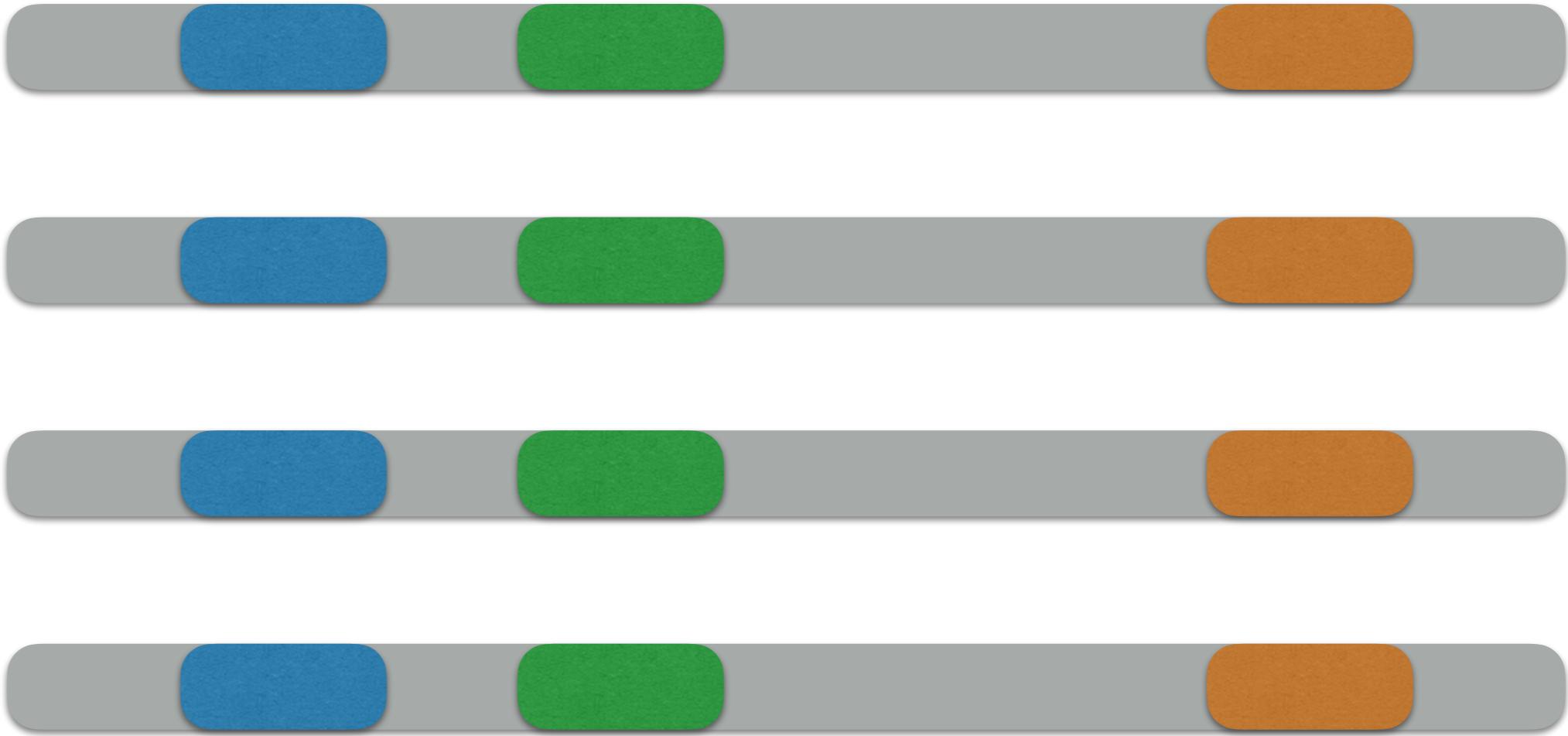


Concatenation or supermatrix



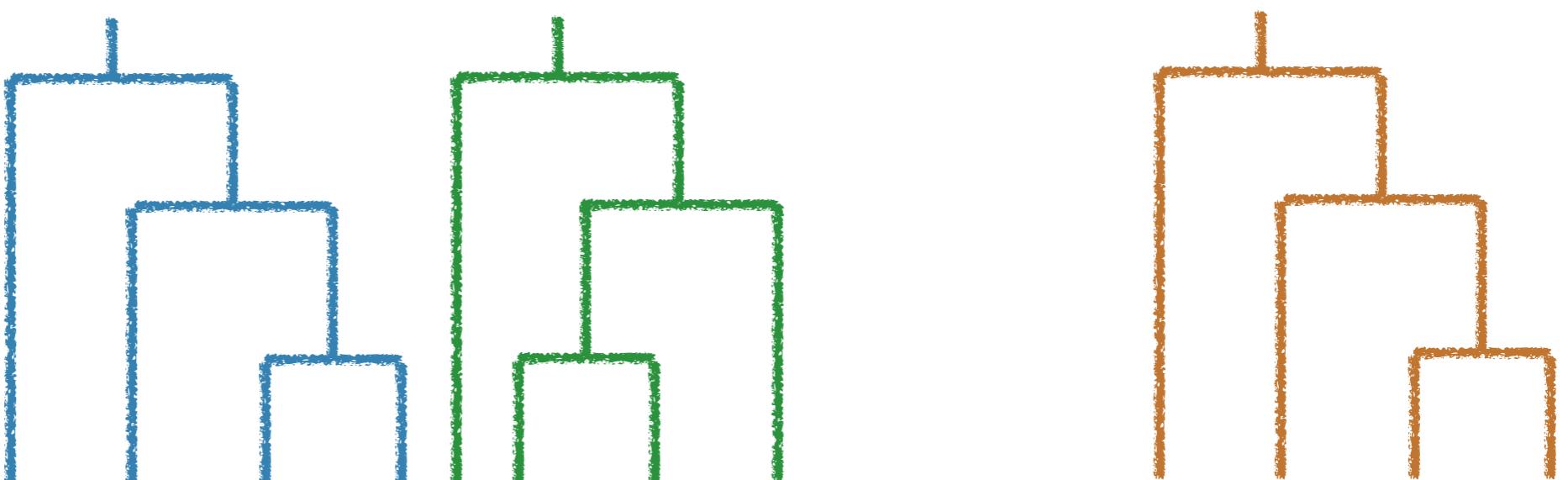


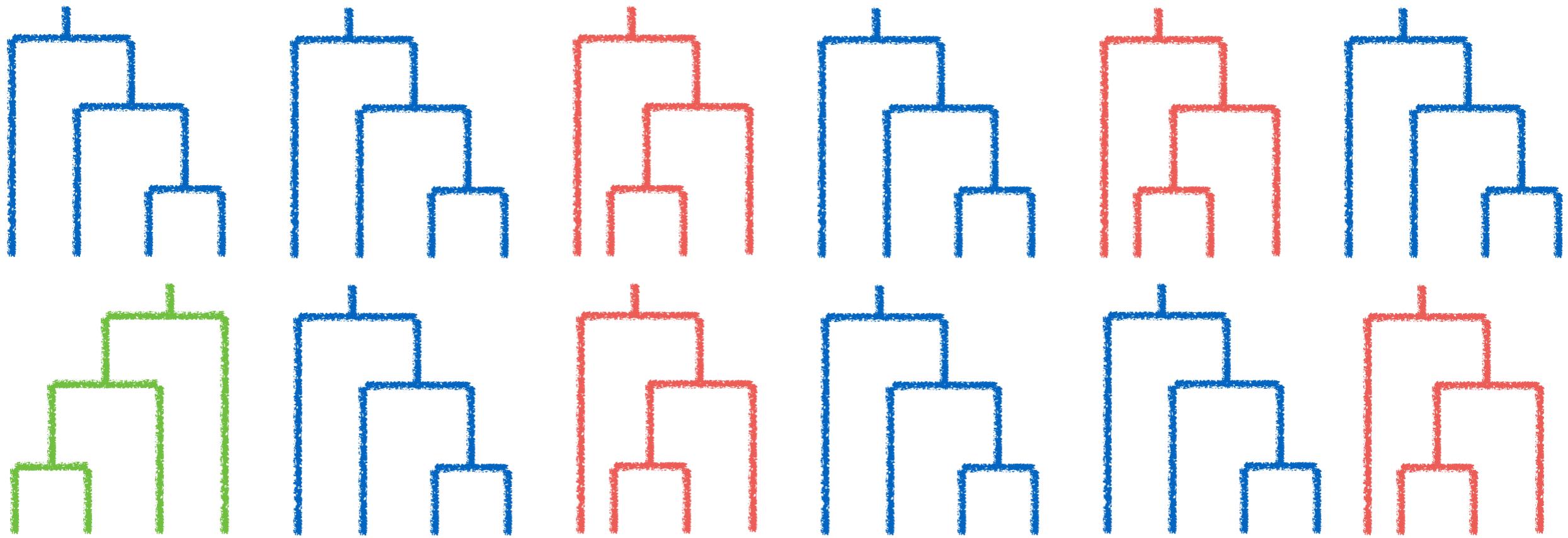
**MrBayes**  
(Huelsenbeck, Ronquist, 2001)  
**RAXML**  
(Stamatakis, 2014)  
**IQ-tree 2**  
(Minh et al, 2020)



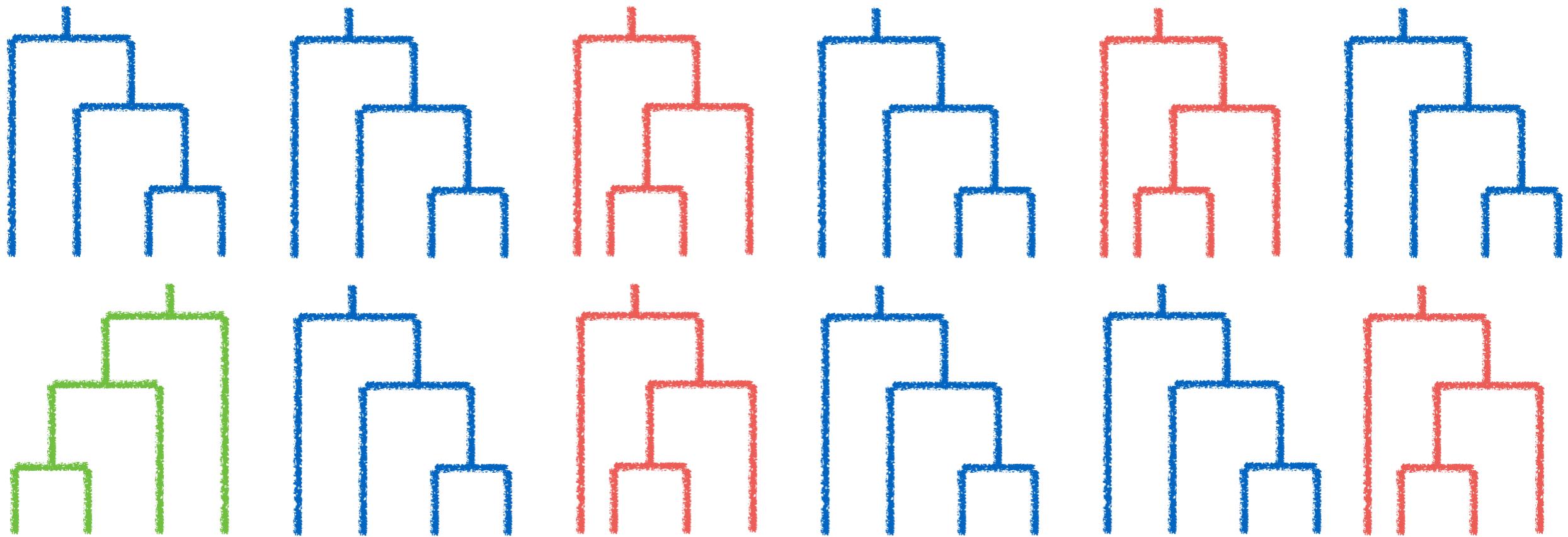
# Estimate gene trees

MrBayes  
(Huelsenbeck, Ronquist, 2001)  
RAxML  
(Stamatakis, 2014)  
IQ-tree 2  
(Minh et al, 2020)



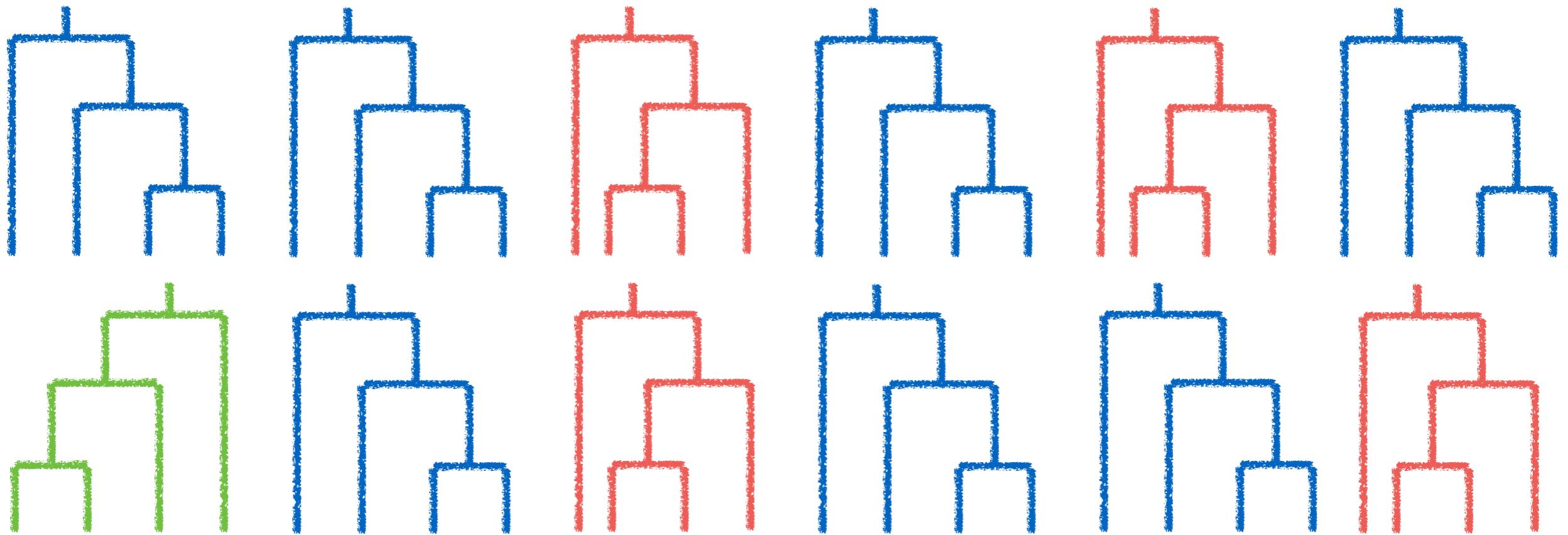


Gene trees



## Gene trees

- Estimation error
- Incomplete lineage sorting
- Gene flow

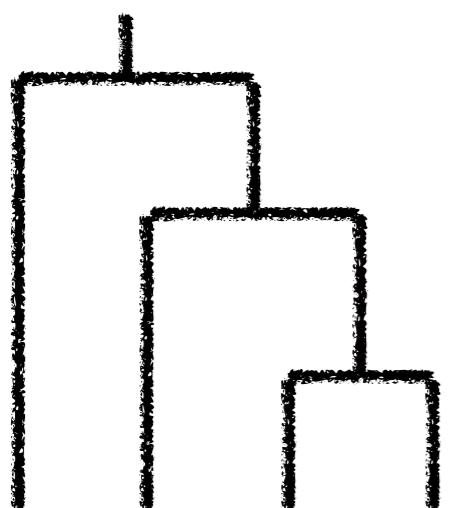
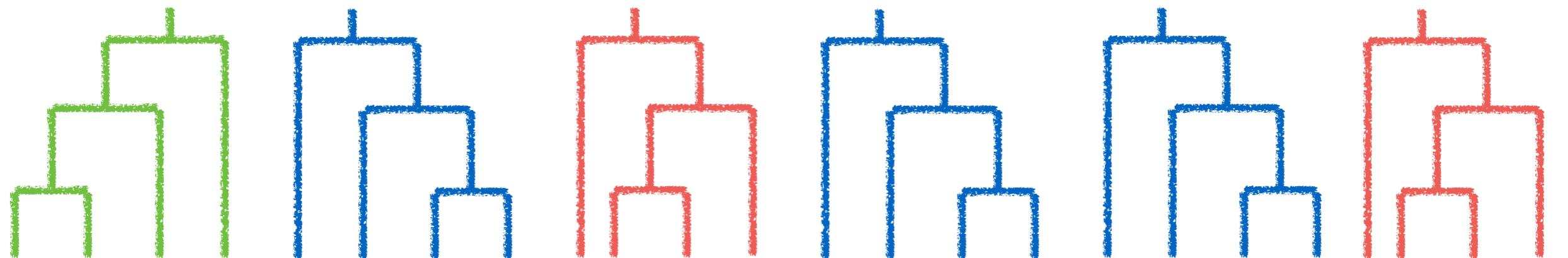
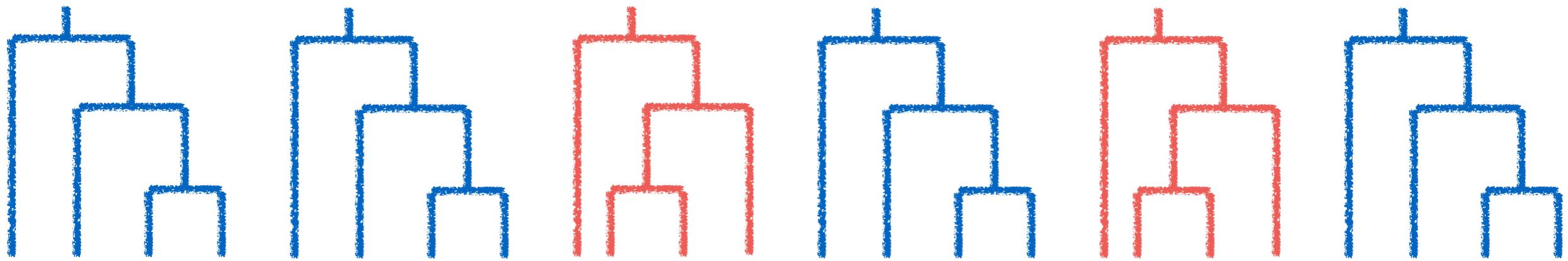


## Gene trees

- Estimation error
- Incomplete lineage sorting
- Gene flow



Species evolutionary  
history



## Gene trees

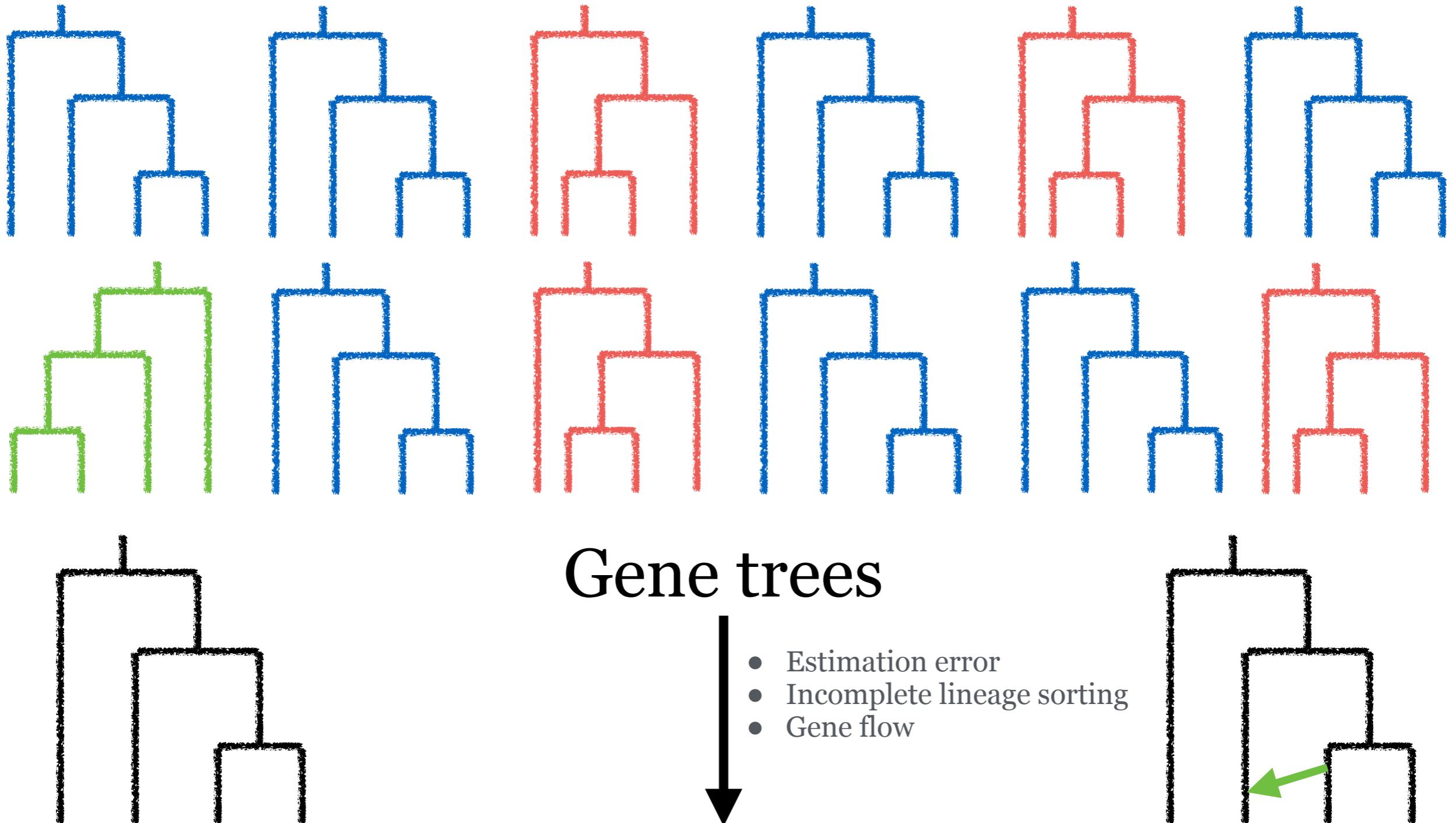


- Estimation error
- Incomplete lineage sorting
- Gene flow

## Species evolutionary history

BUCKY (Ané et al, 2007)

ASTRAL (Mirarab et al, 2014)

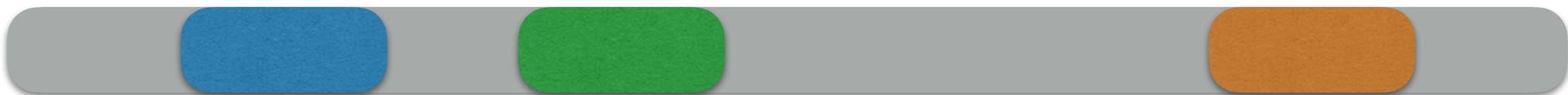


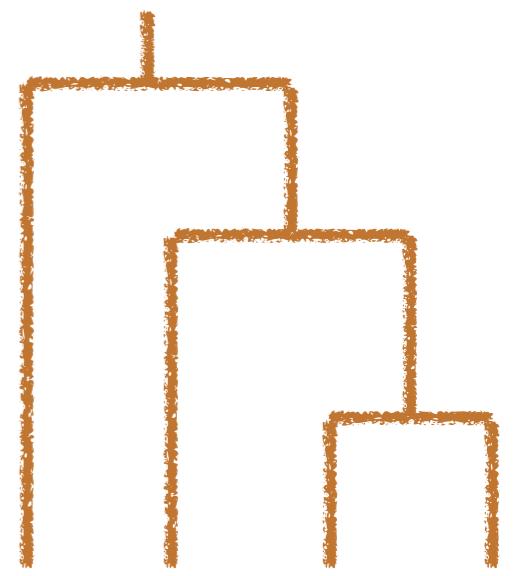
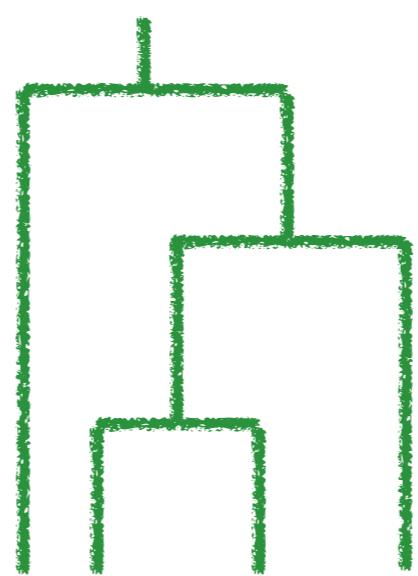
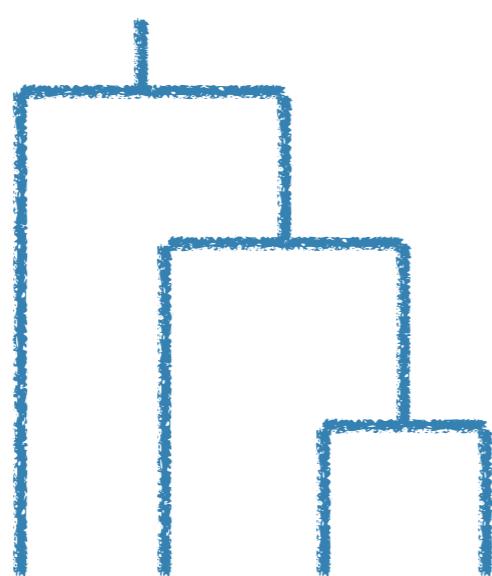
Species evolutionary  
history



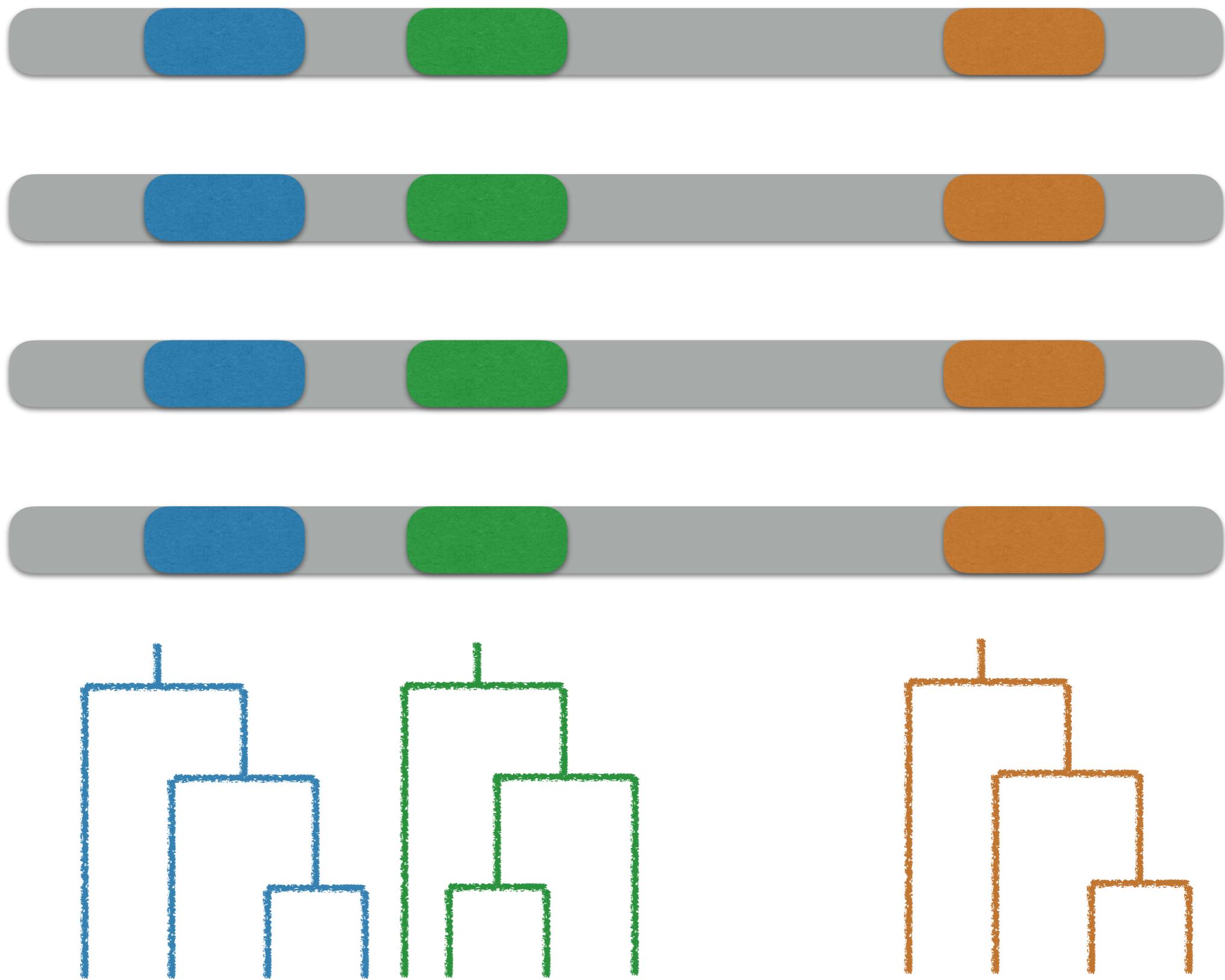
BUCKY (Ané et al, 2007)  
ASTRAL (Mirarab et al, 2014)

SNaQ  
(Solís-Lemus, Ané, 2016)

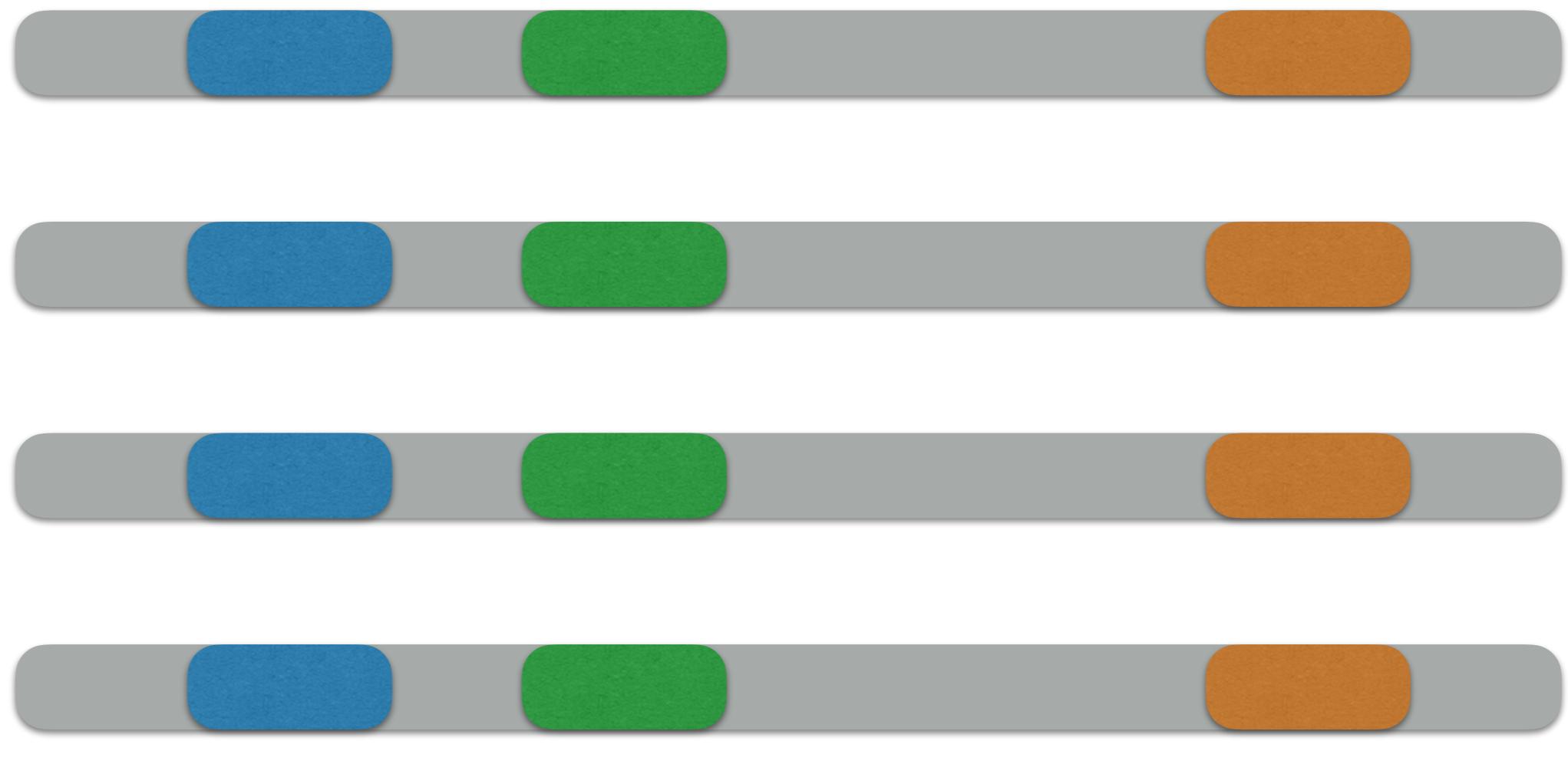




Ortholog  
Recombination-free  
MDL (Ané, 2011, GBE)



Ortholog  
Recombination-free  
MDL (Ané, 2011, GBE)

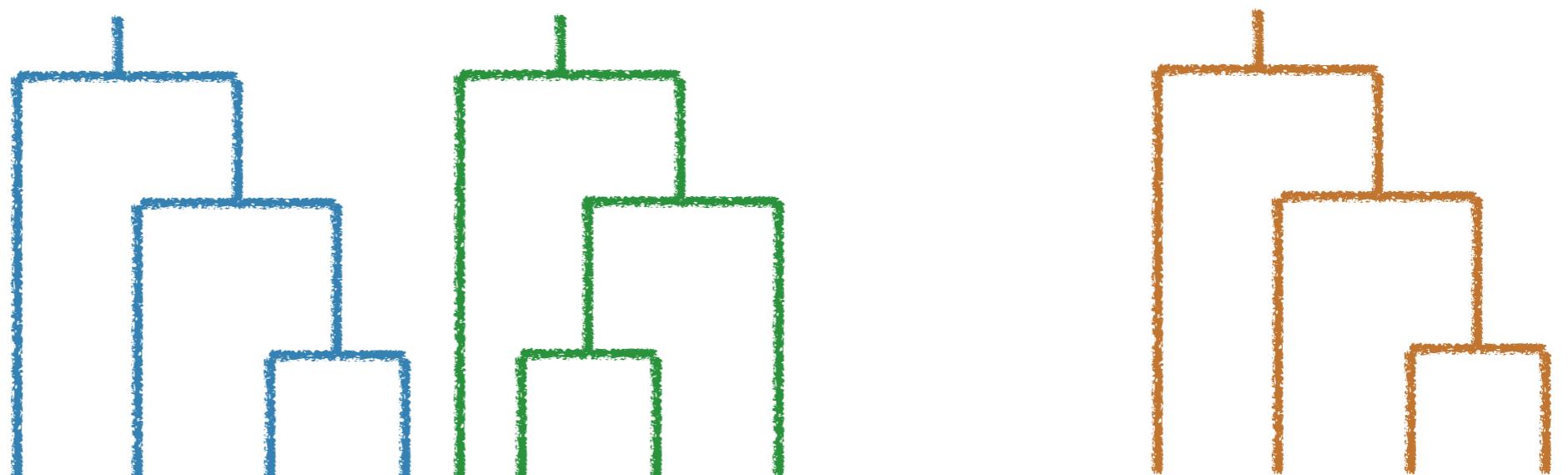


# Estimate gene trees

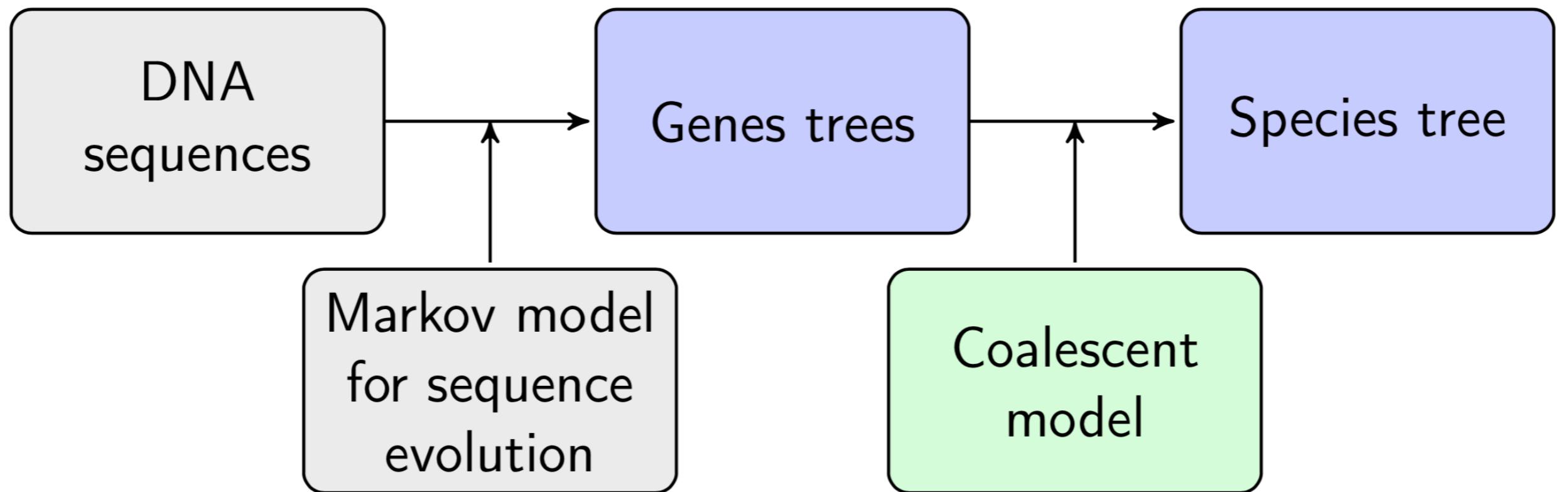
MrBayes  
(Huelsenbeck, Ronquist, 2001)

RAXML  
(Stamatakis, 2014)

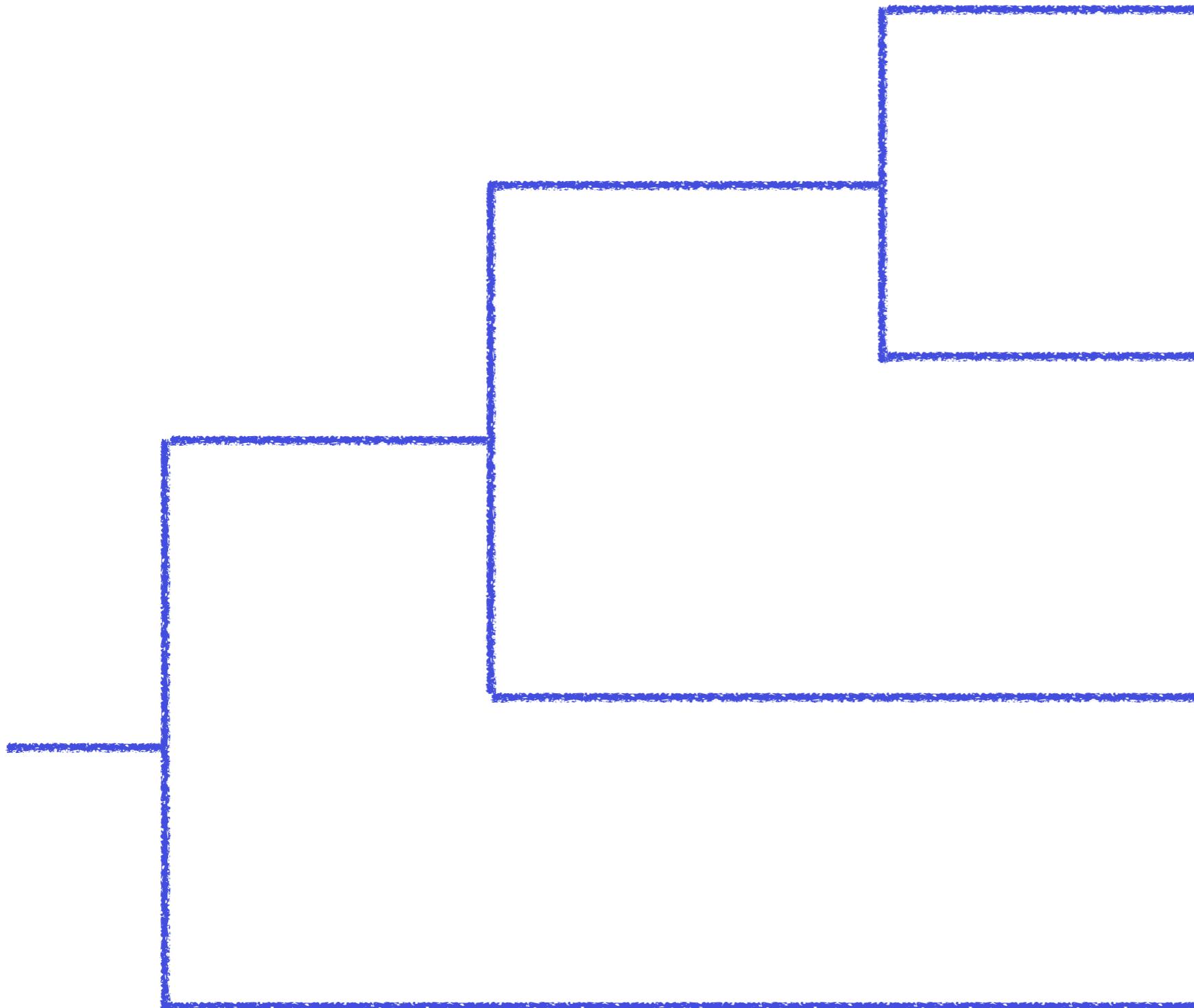
IQ-tree 2  
(Minh et al, 2020)



# Phylogenetic inference

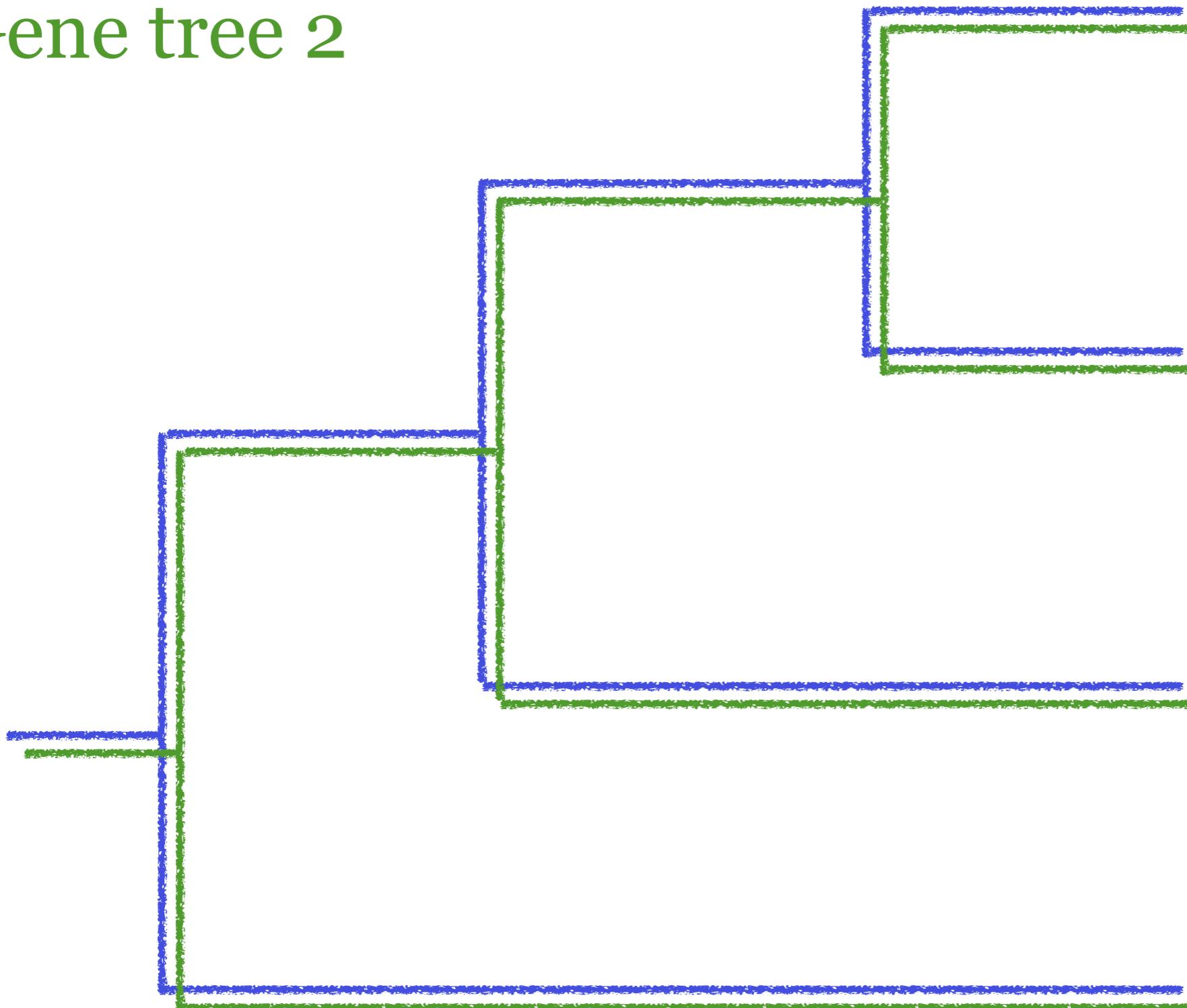


# Gene tree



Gene tree 1

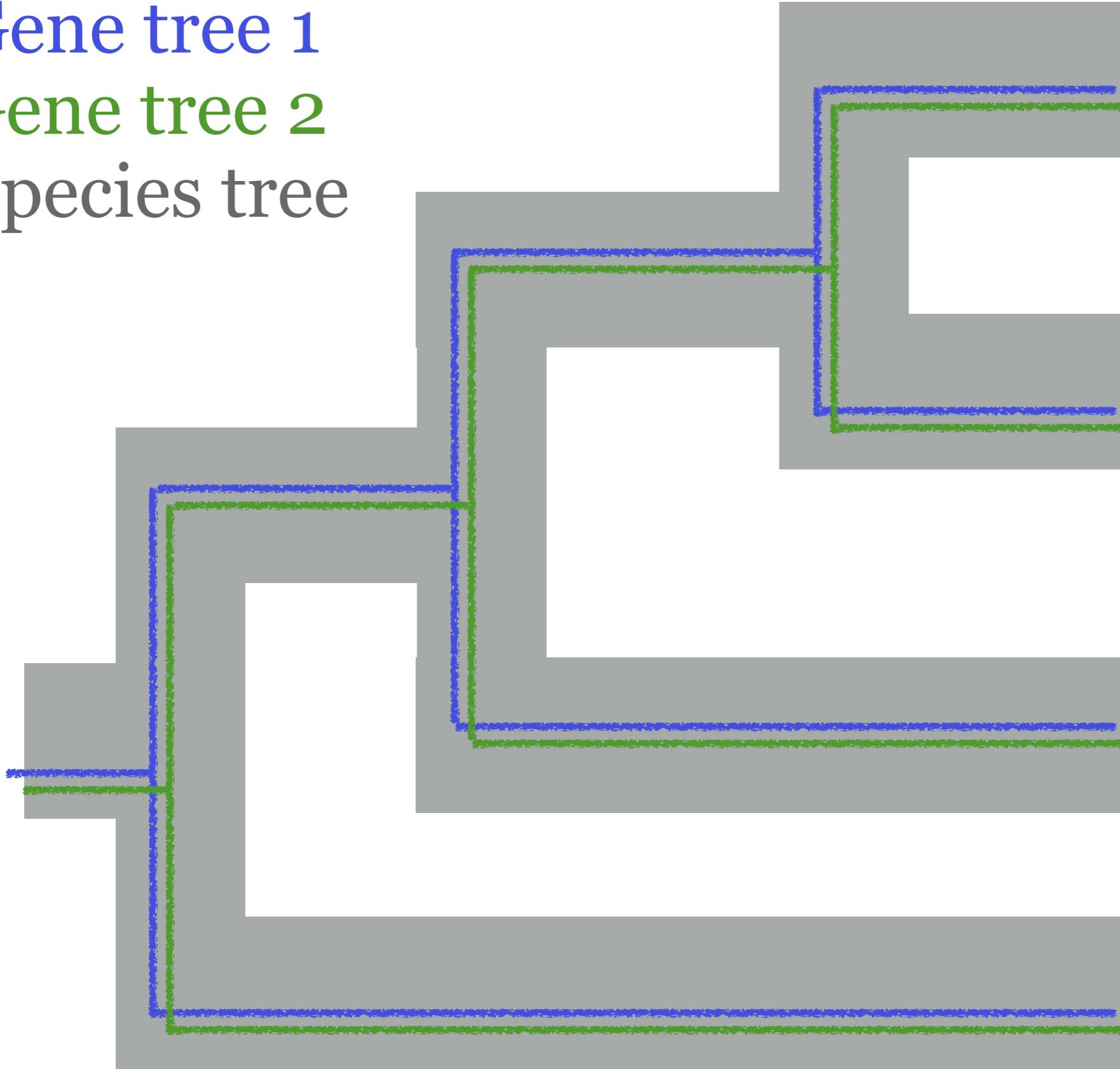
Gene tree 2



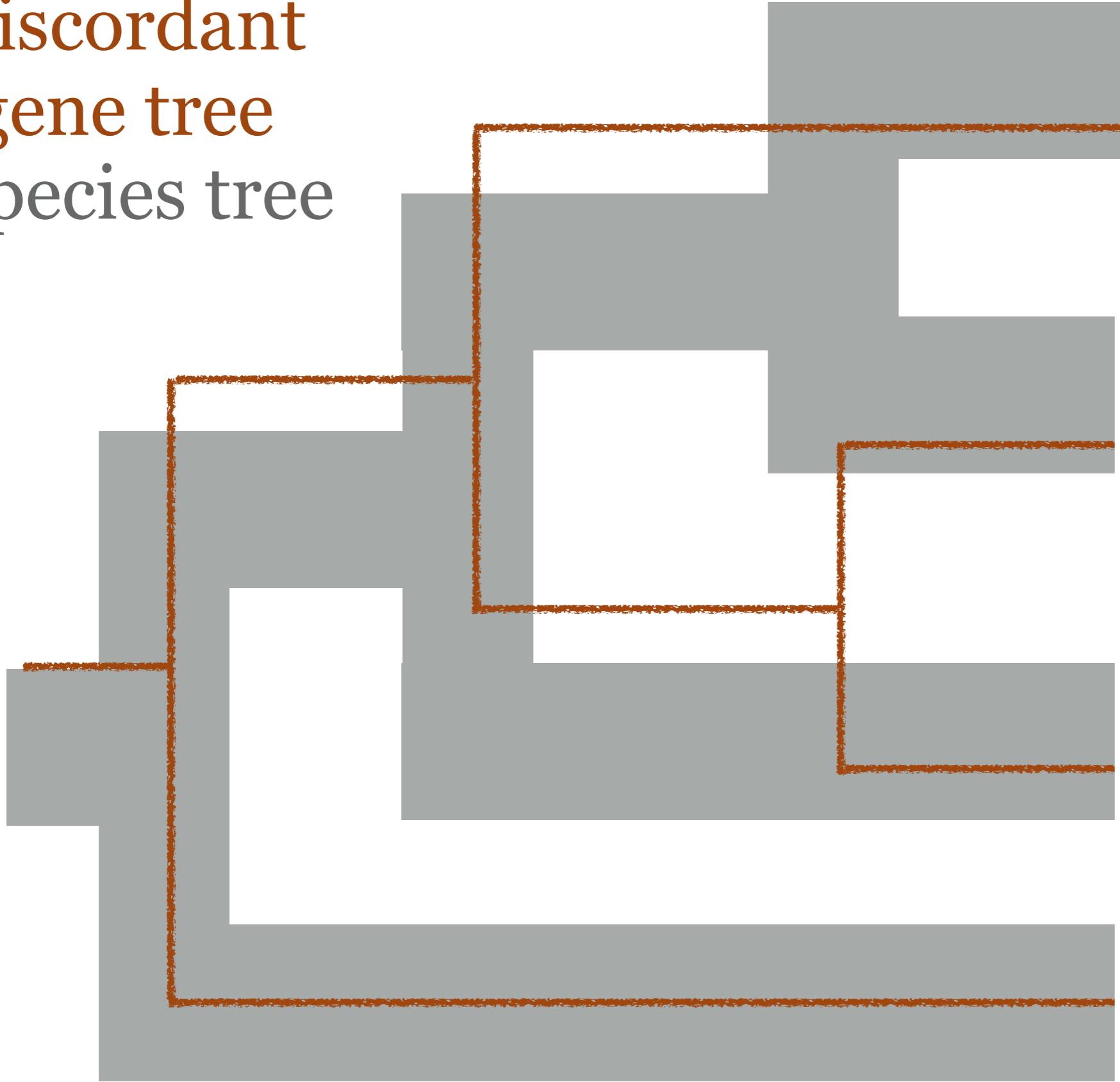
# Gene tree 1

# Gene tree 2

# Species tree

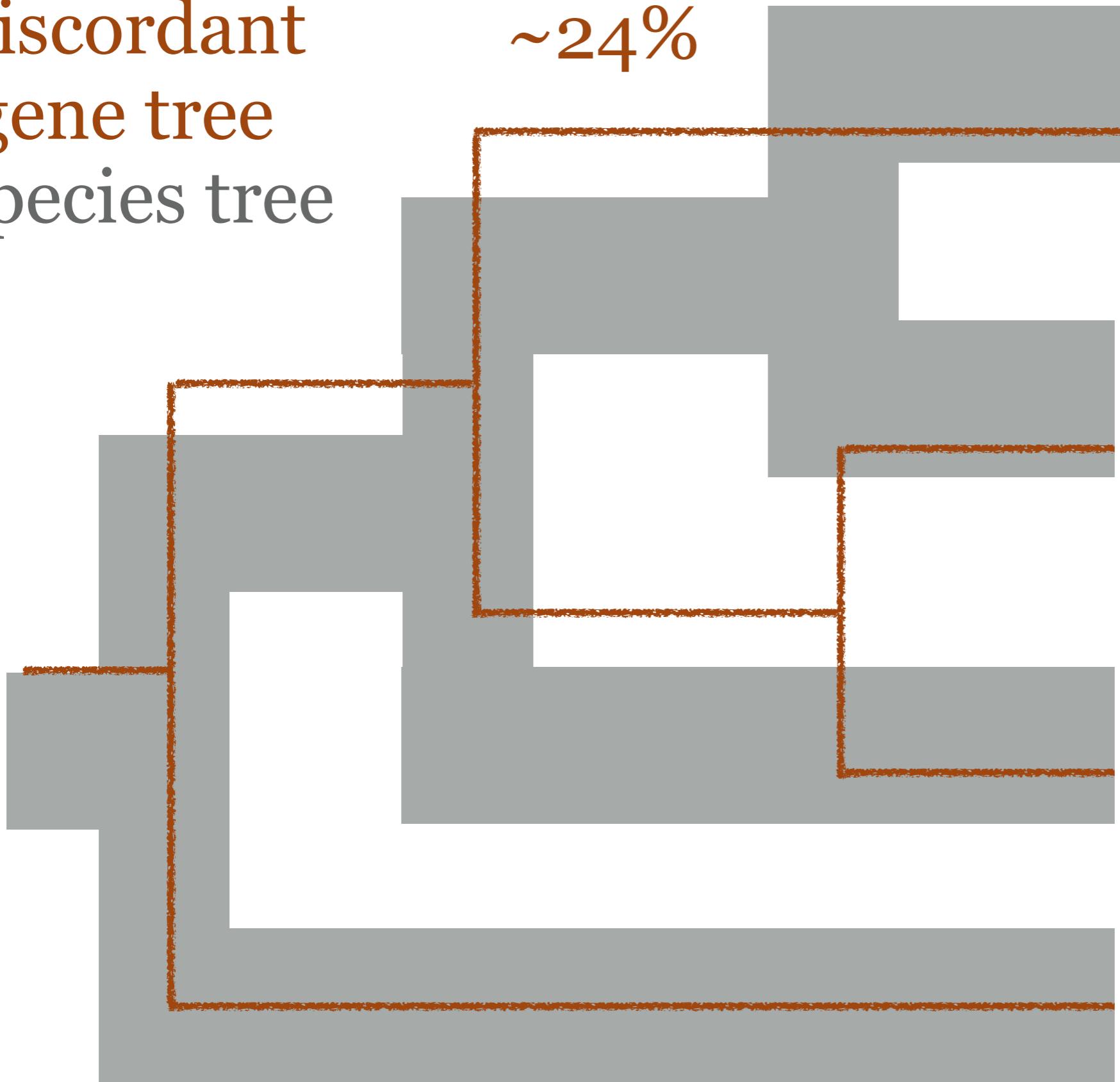


# Discordant gene tree Species tree



# Discordant gene tree Species tree

~24%



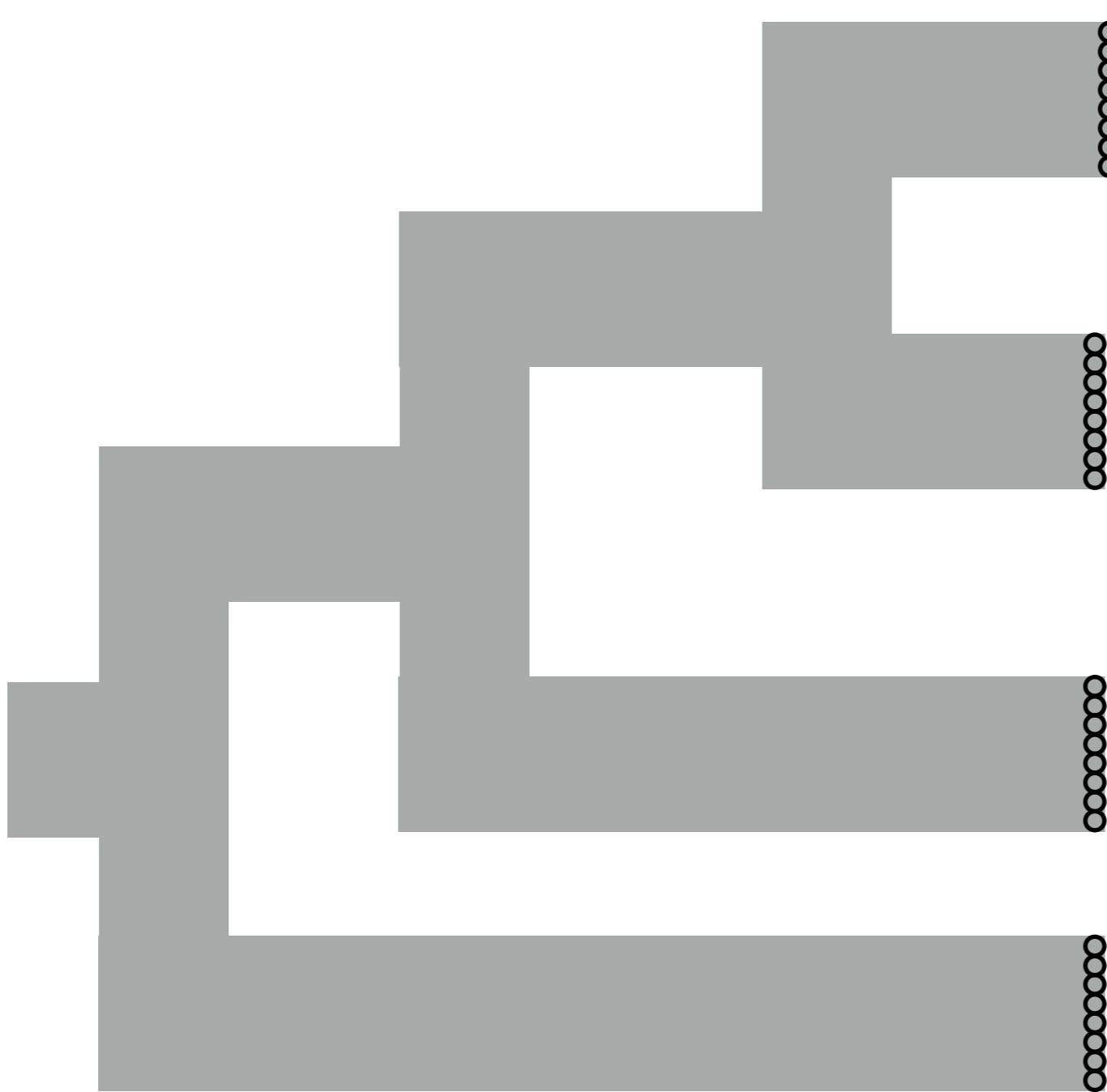
# Coalescent model



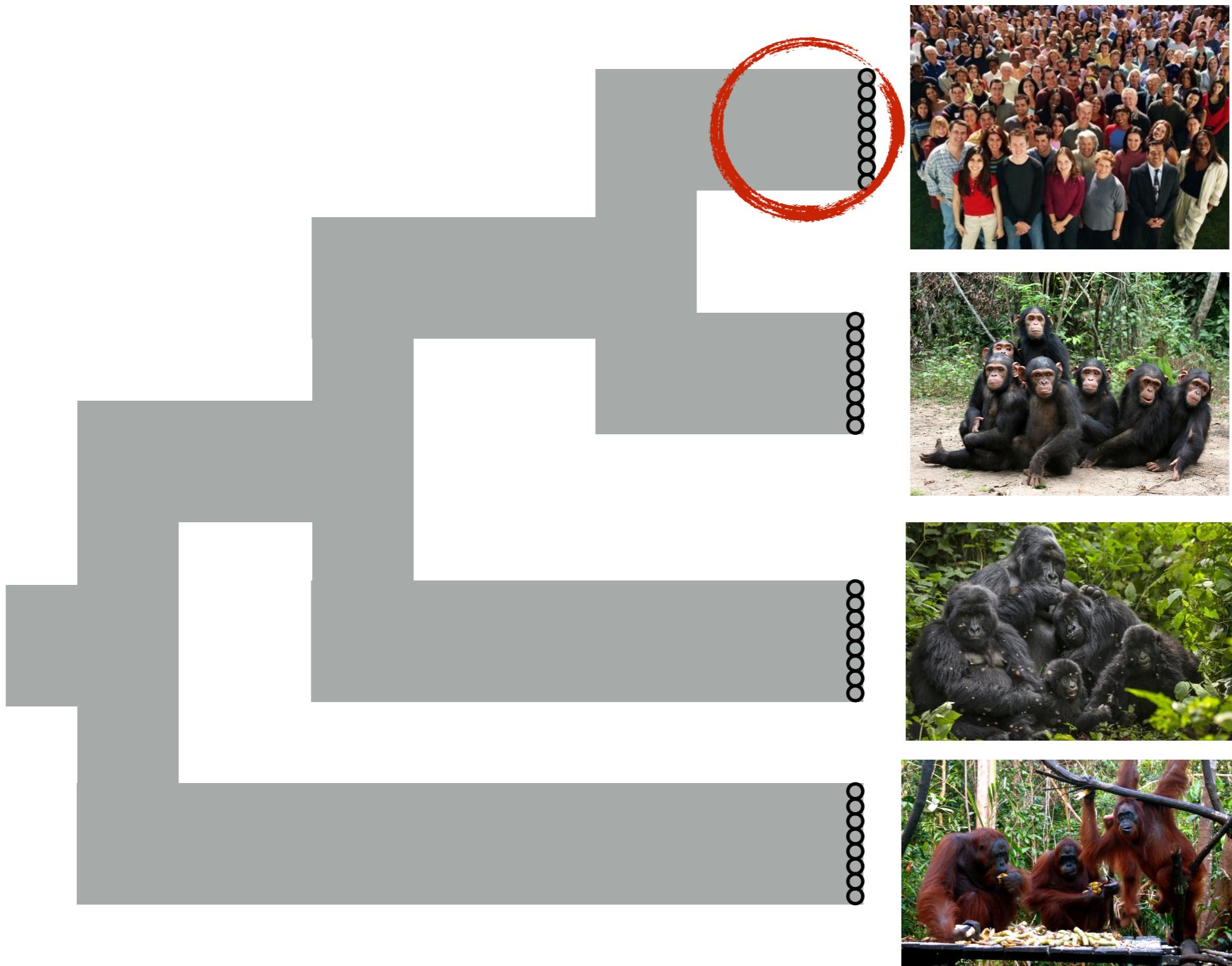
# Coalescent model



# Coalescent model

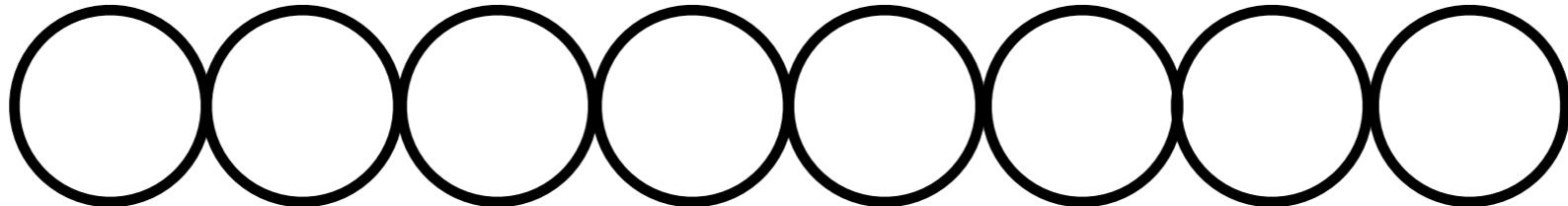


# Coalescent model



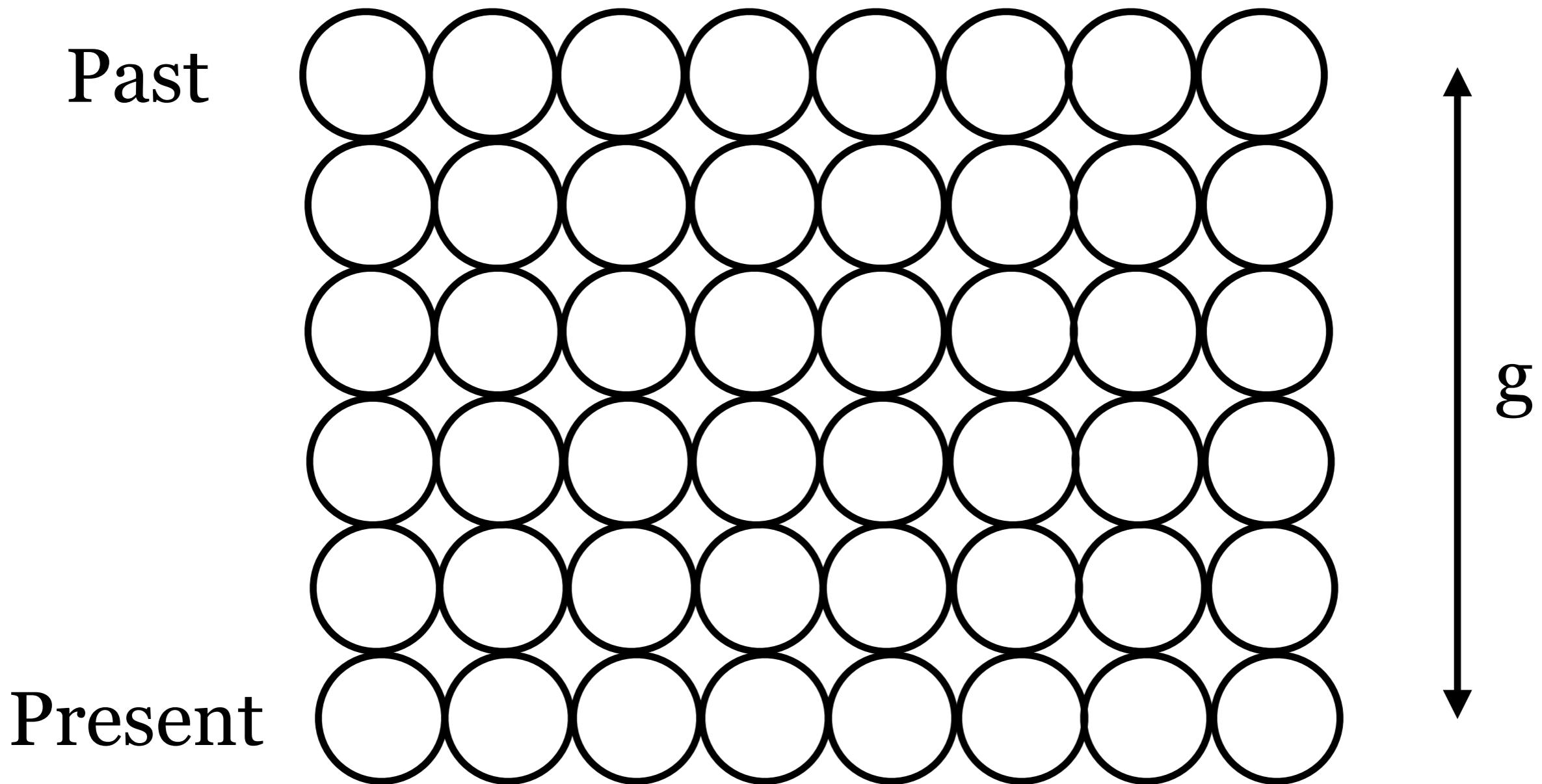
# Coalescent model within I population

Present



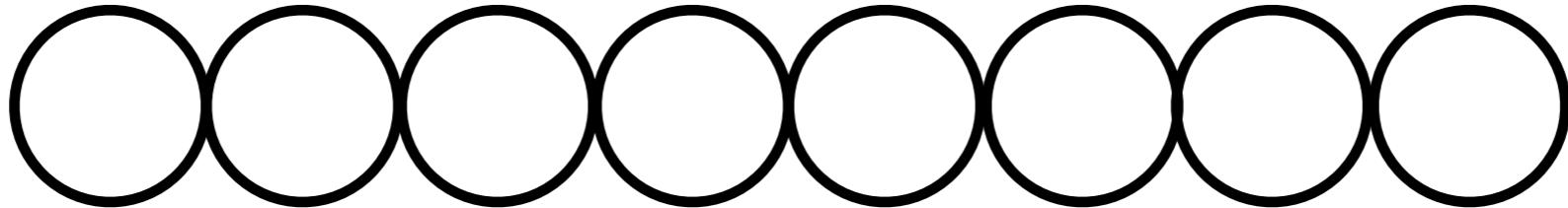
# Coalescent model within I

population



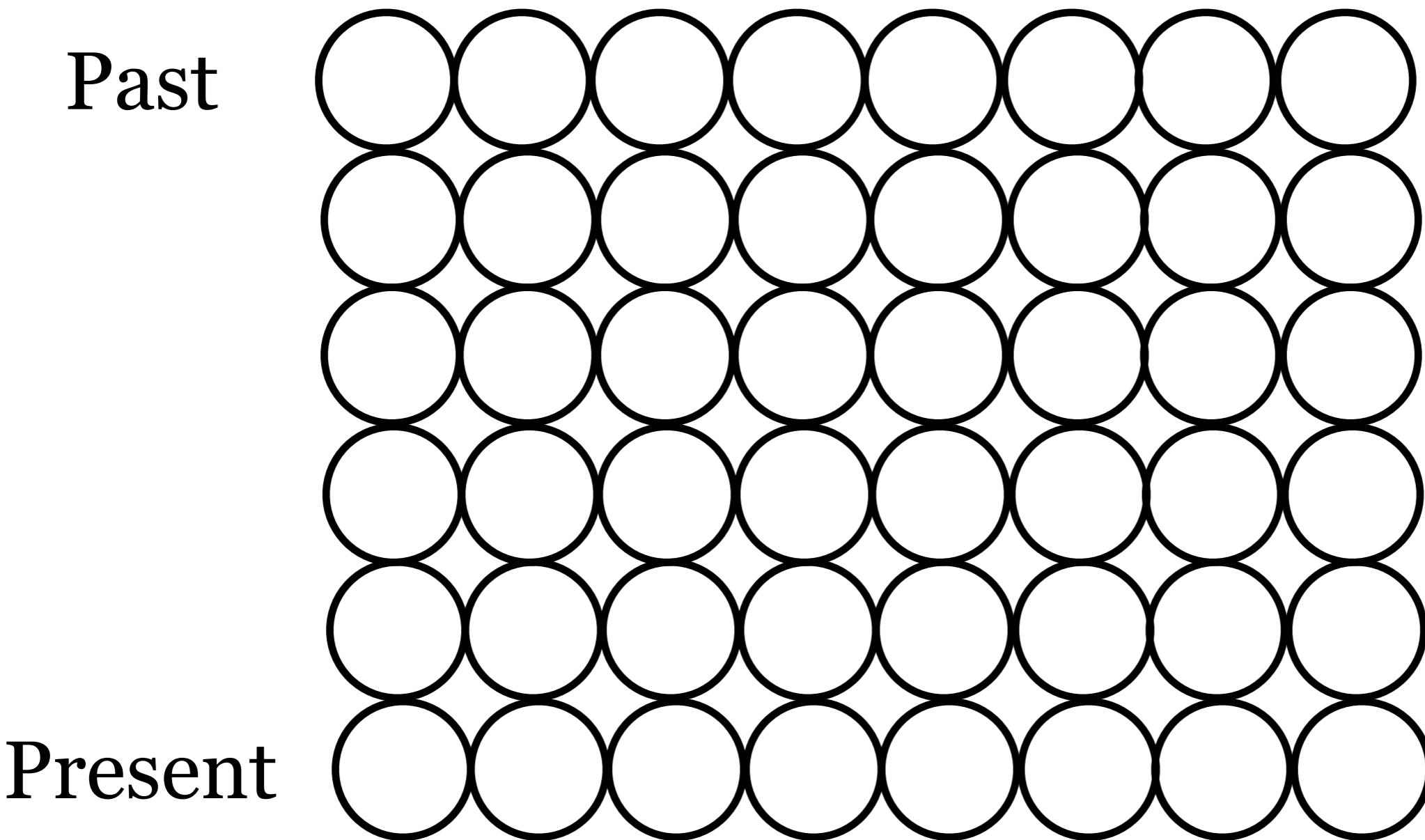
# Coalescent model within I population

Present



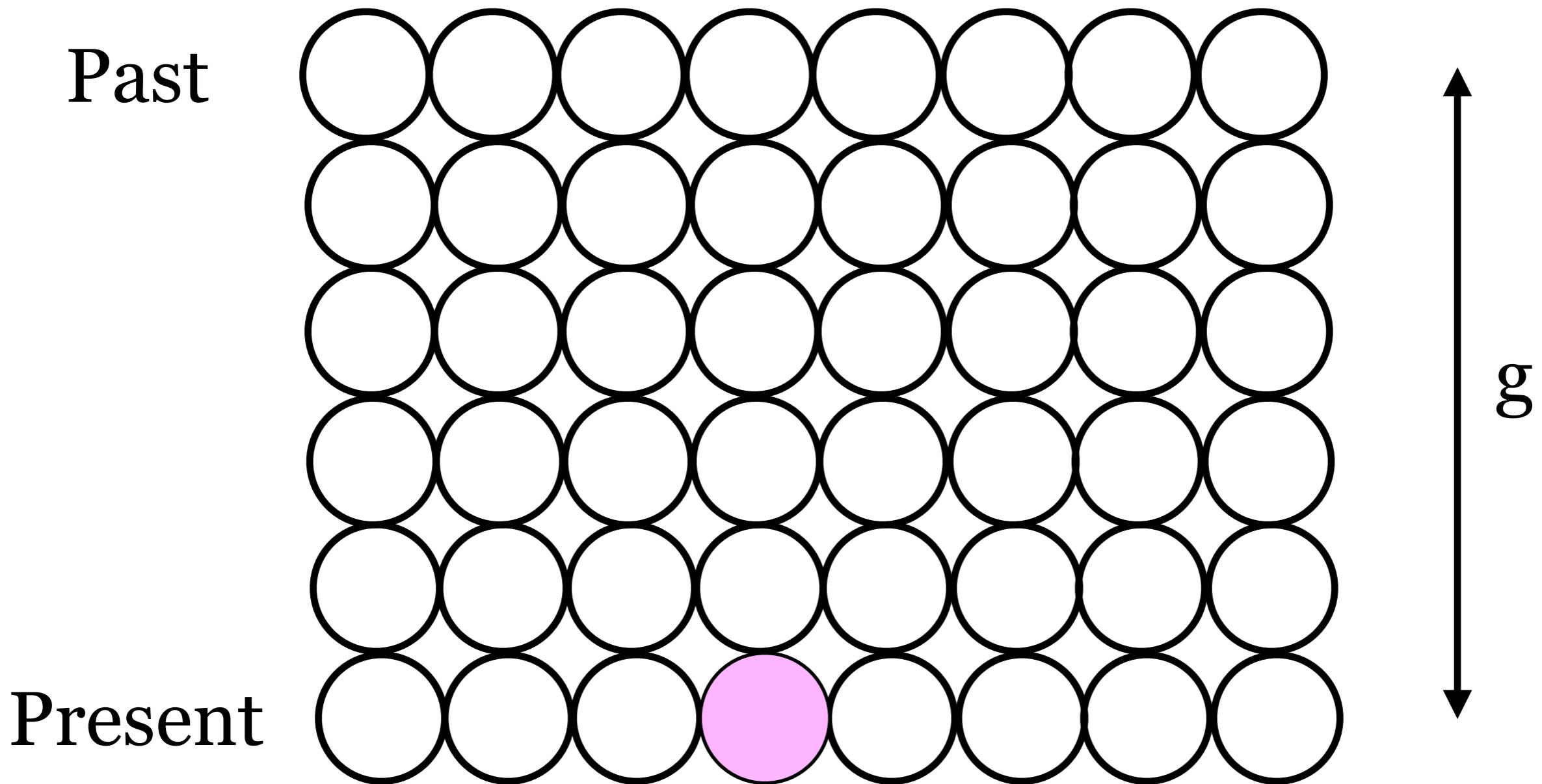
# Coalescent model within I

population



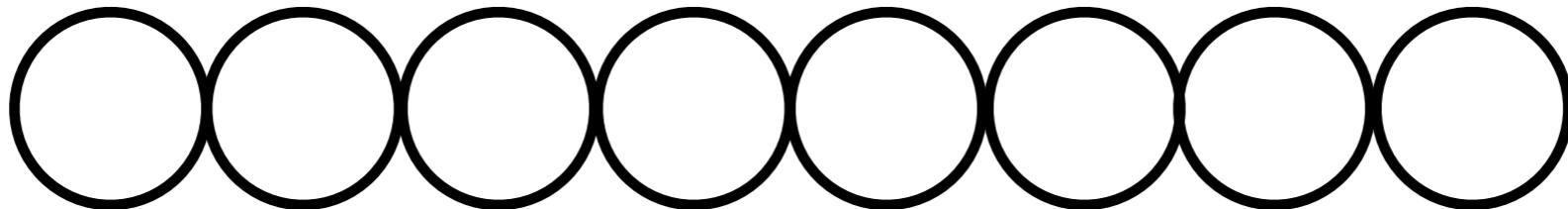
# Coalescent model within I

population



# Coalescent model within I population

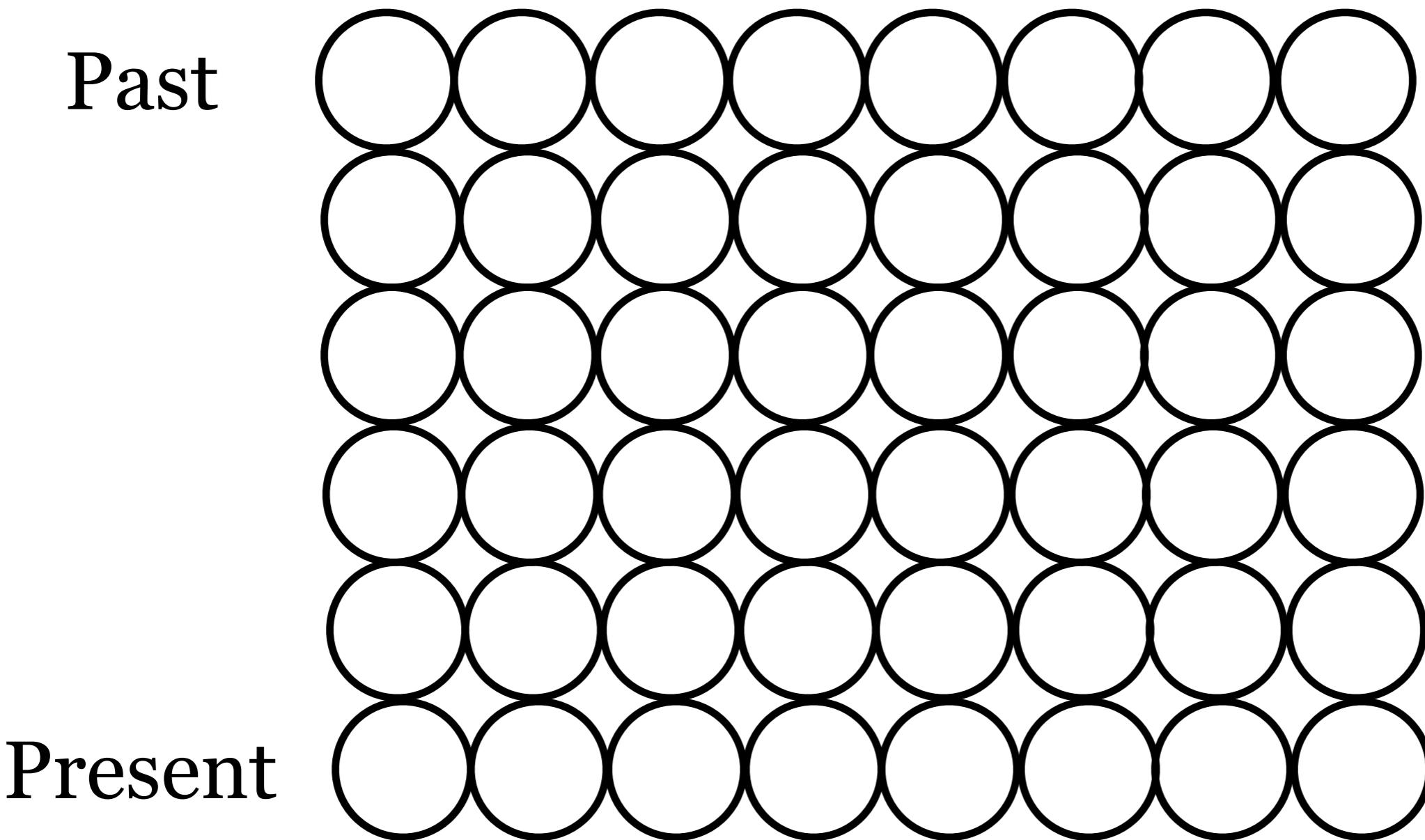
Present



No selection: uniform probability  $1/N$

# Coalescent model within I

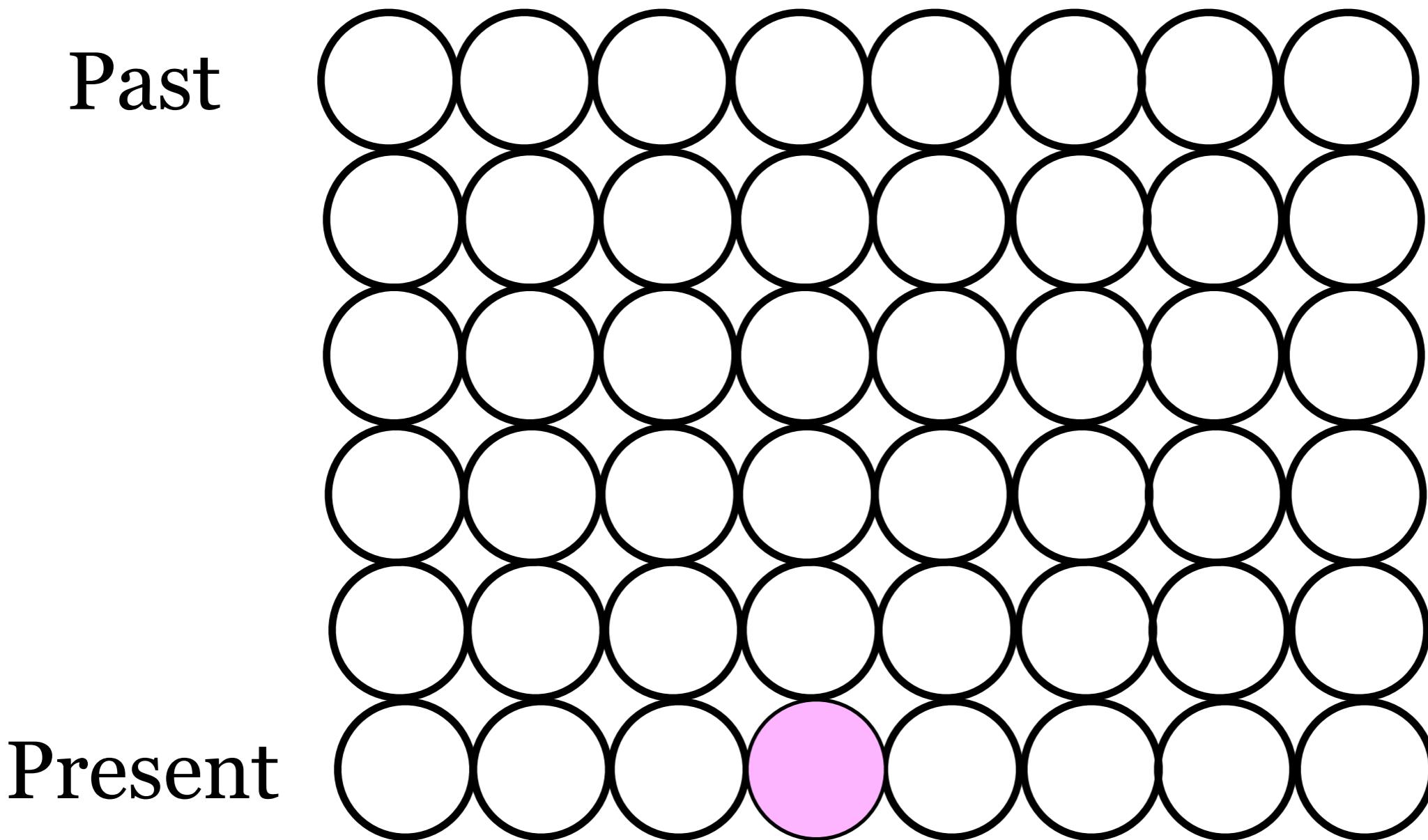
population



No selection: uniform probability  $1/N$

# Coalescent model within I

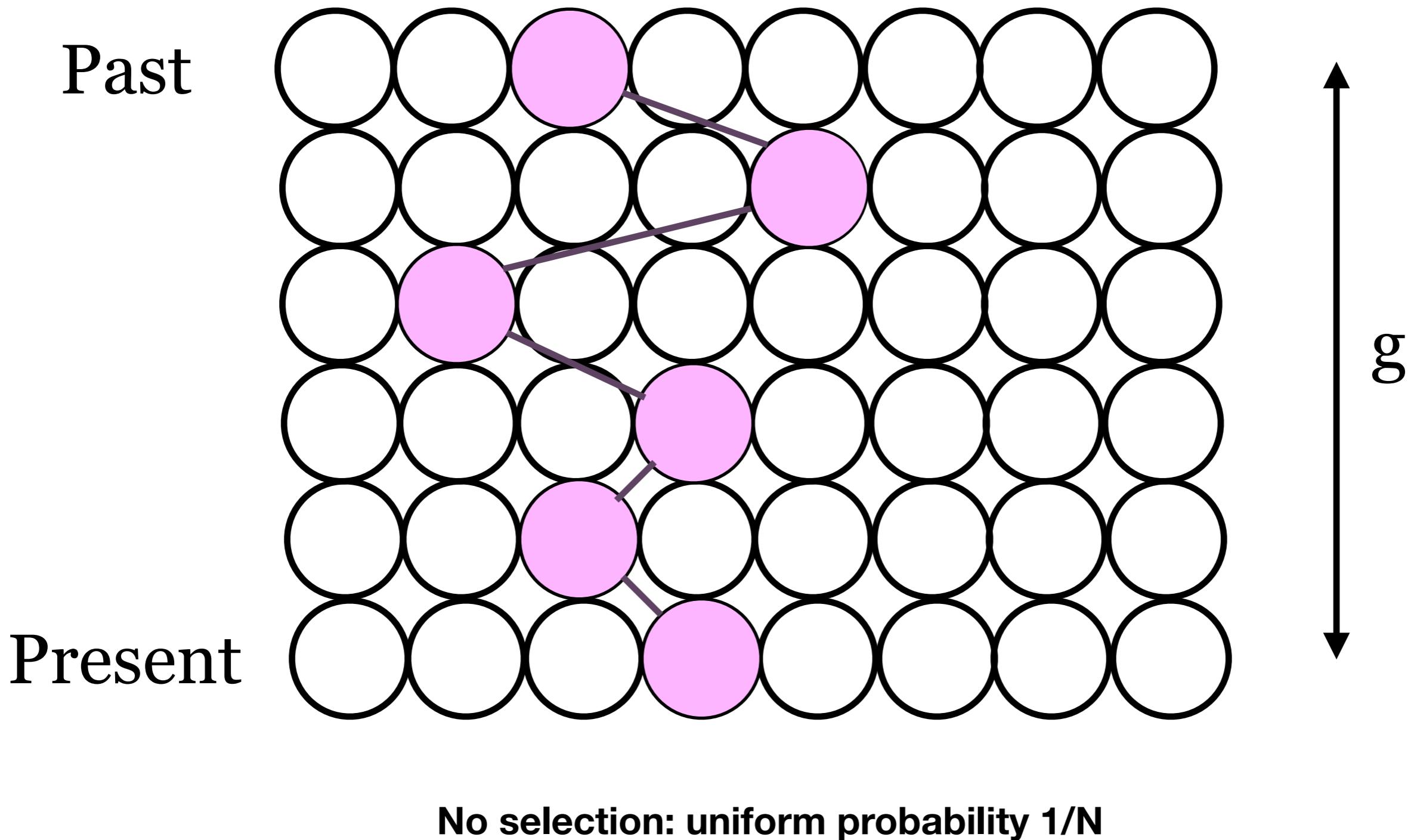
population



No selection: uniform probability  $1/N$

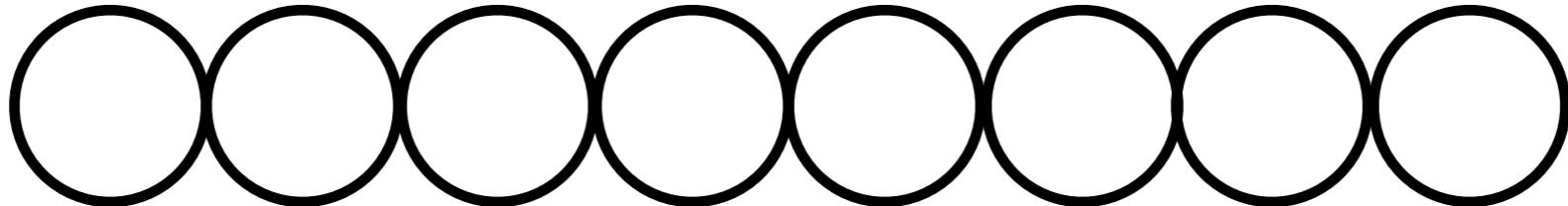
# Coalescent model within I

population



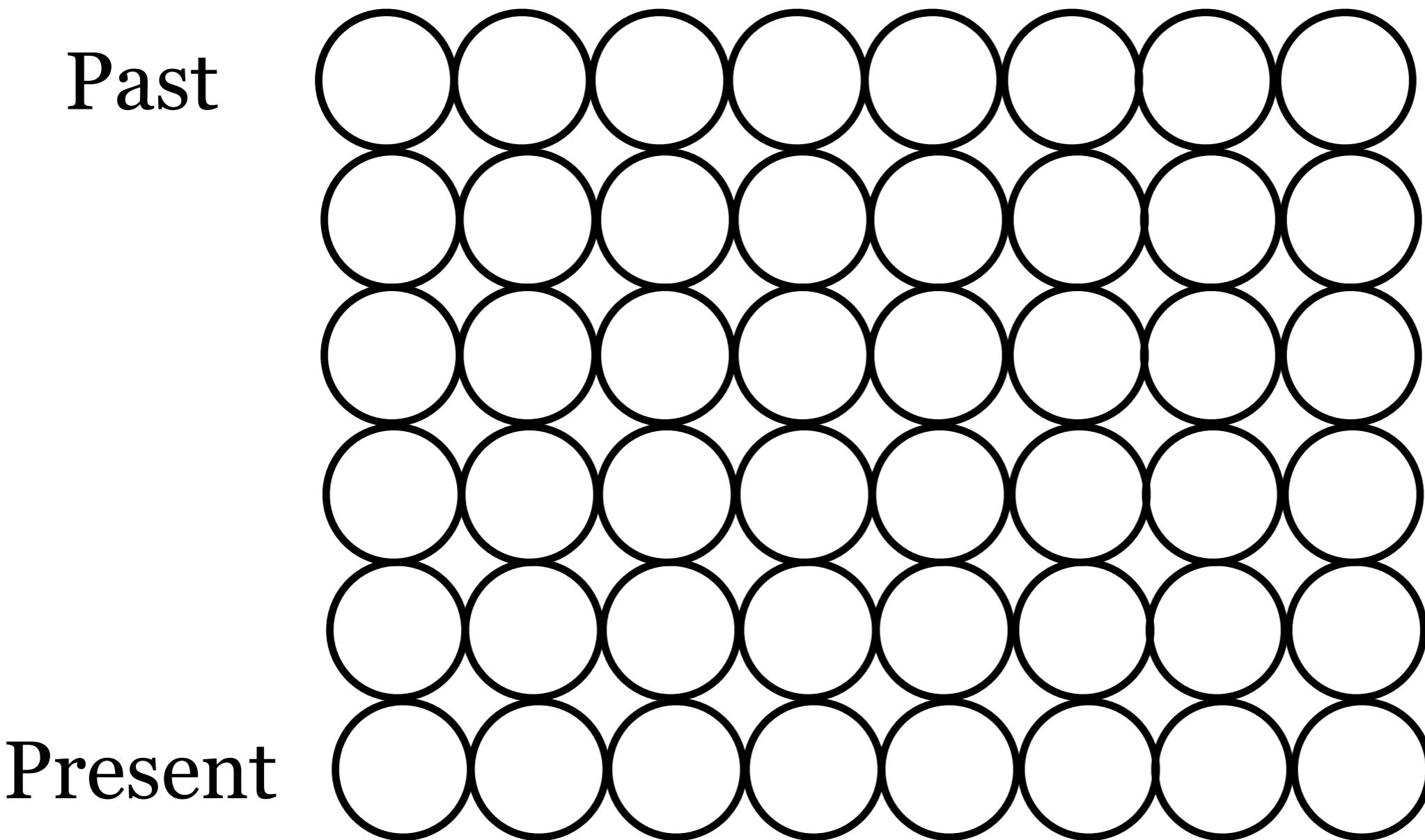
# Coalescent model within I population

Present



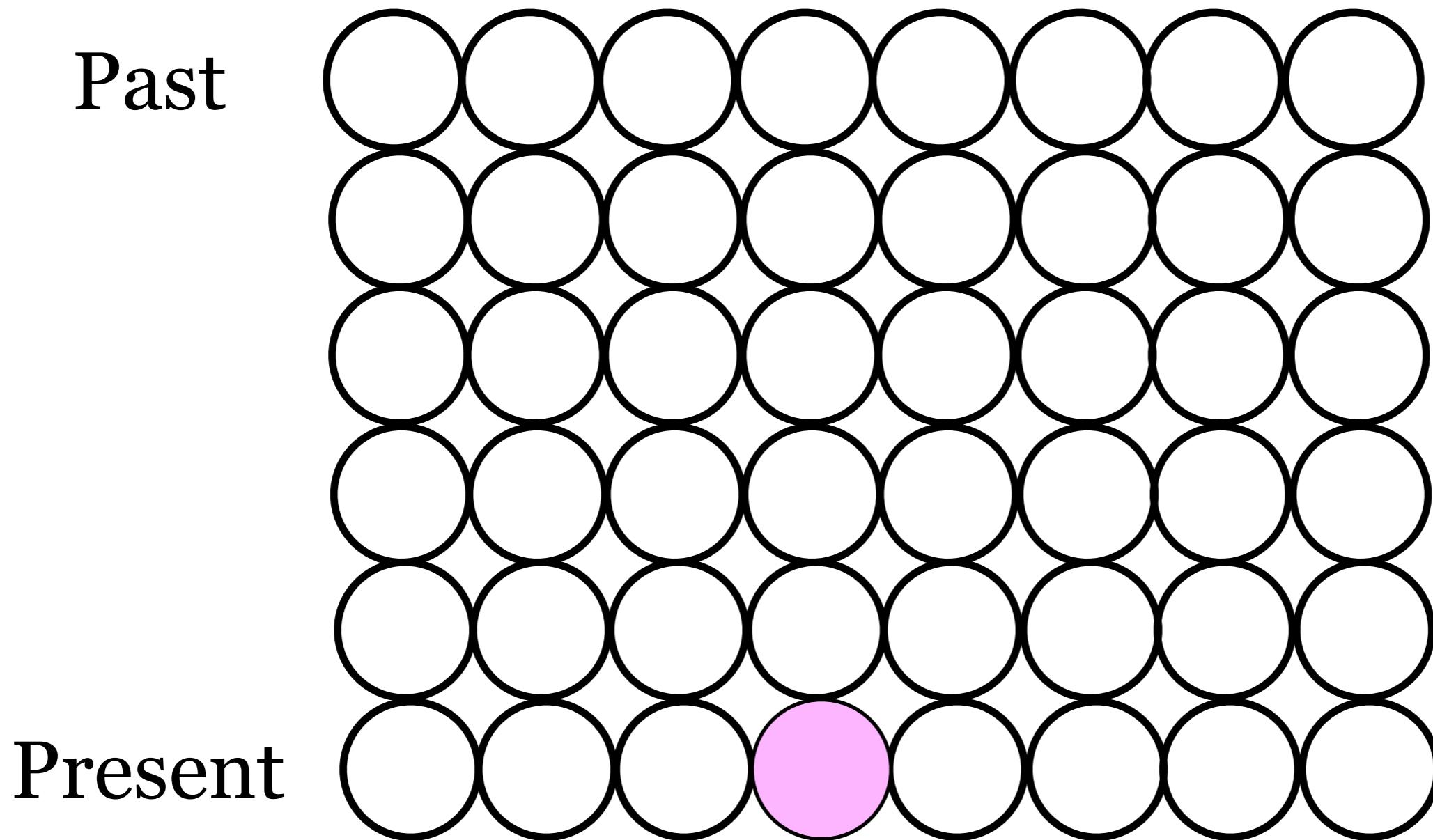
# Coalescent model within I

population



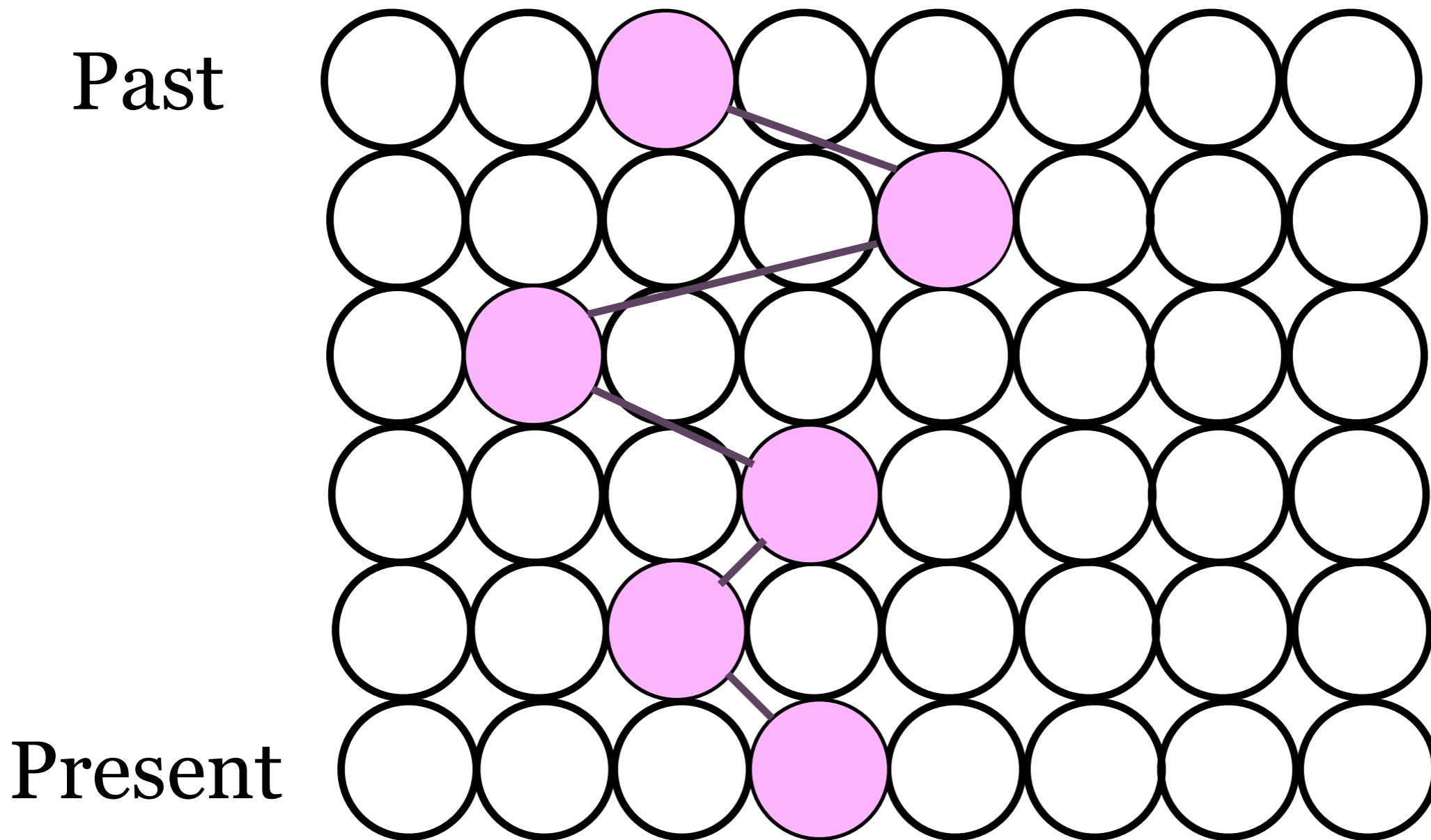
# Coalescent model within I

population



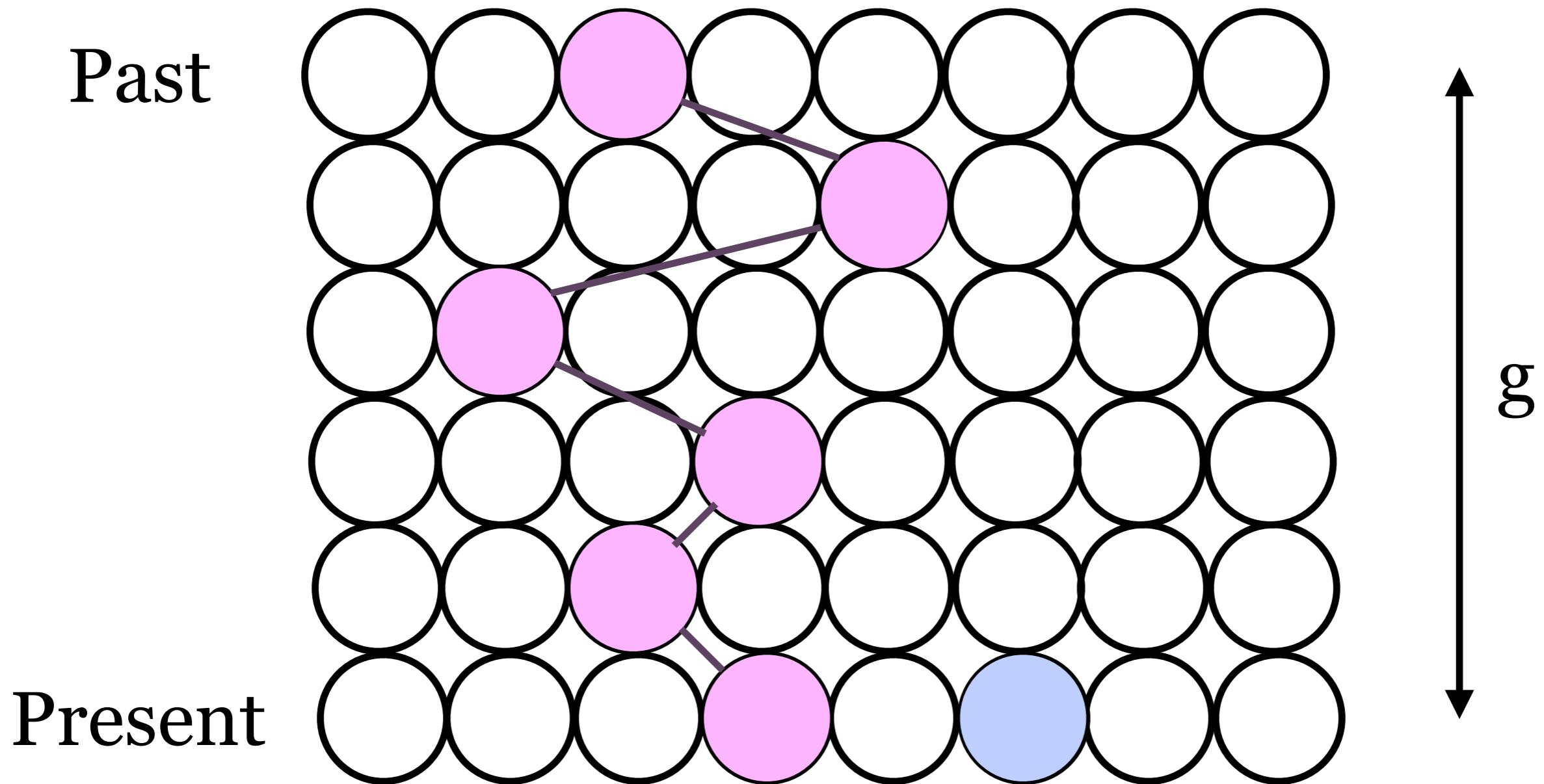
# Coalescent model within I

population



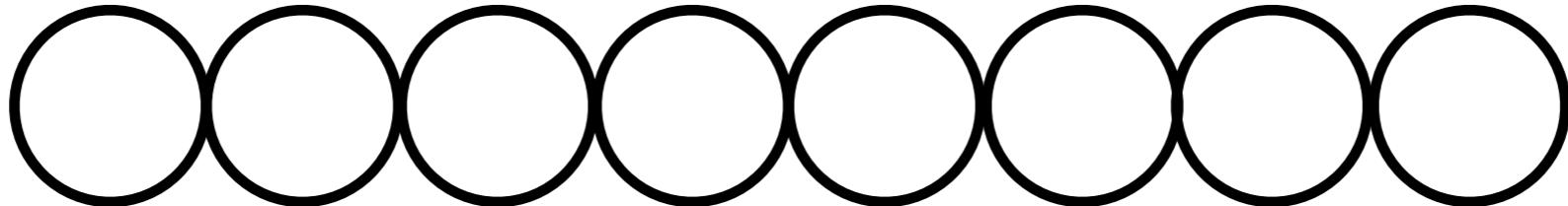
# Coalescent model within I

population



# Coalescent model within I population

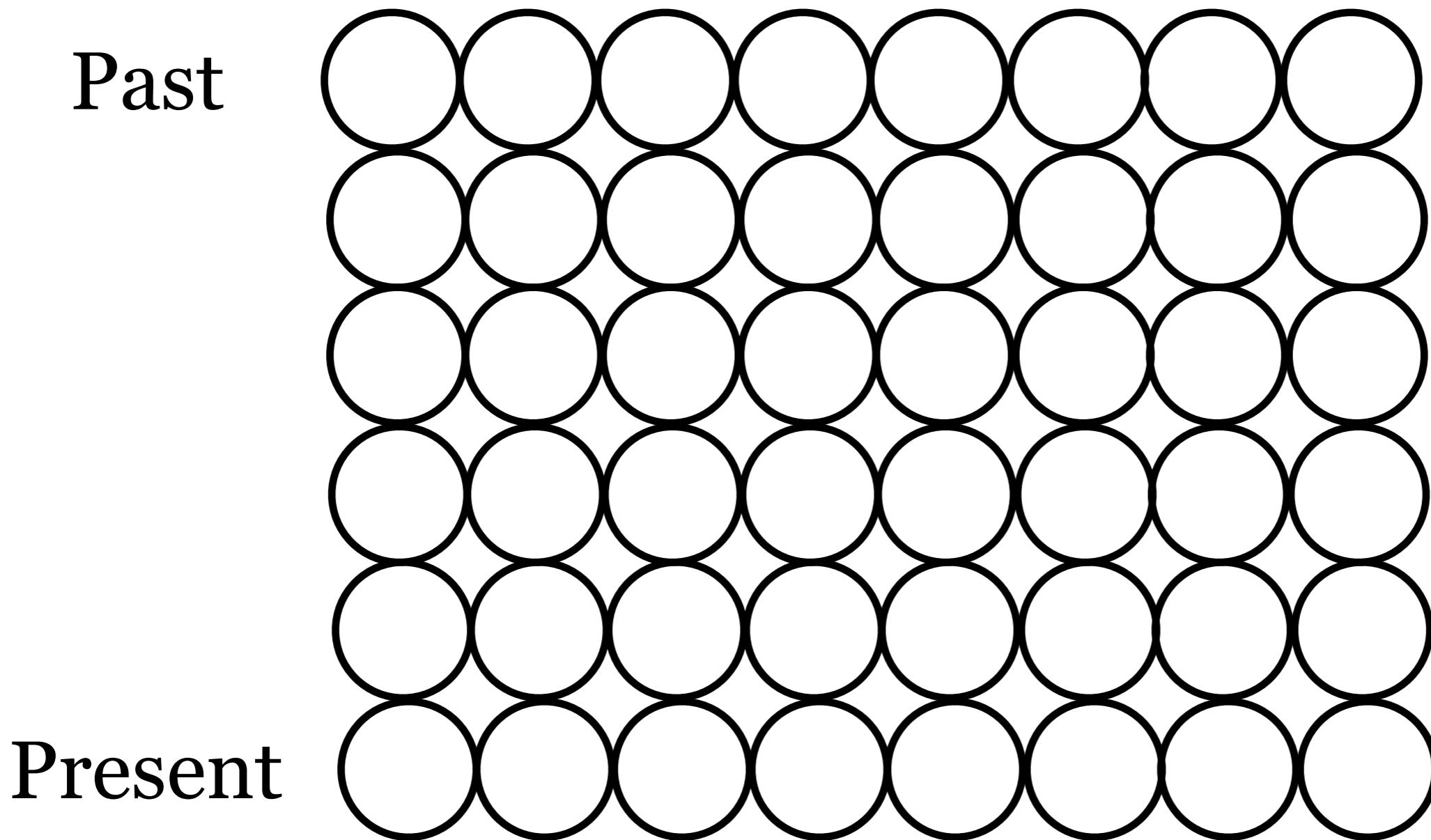
Present



**How many generations do we have to wait for these two individuals to reach a common ancestor?**

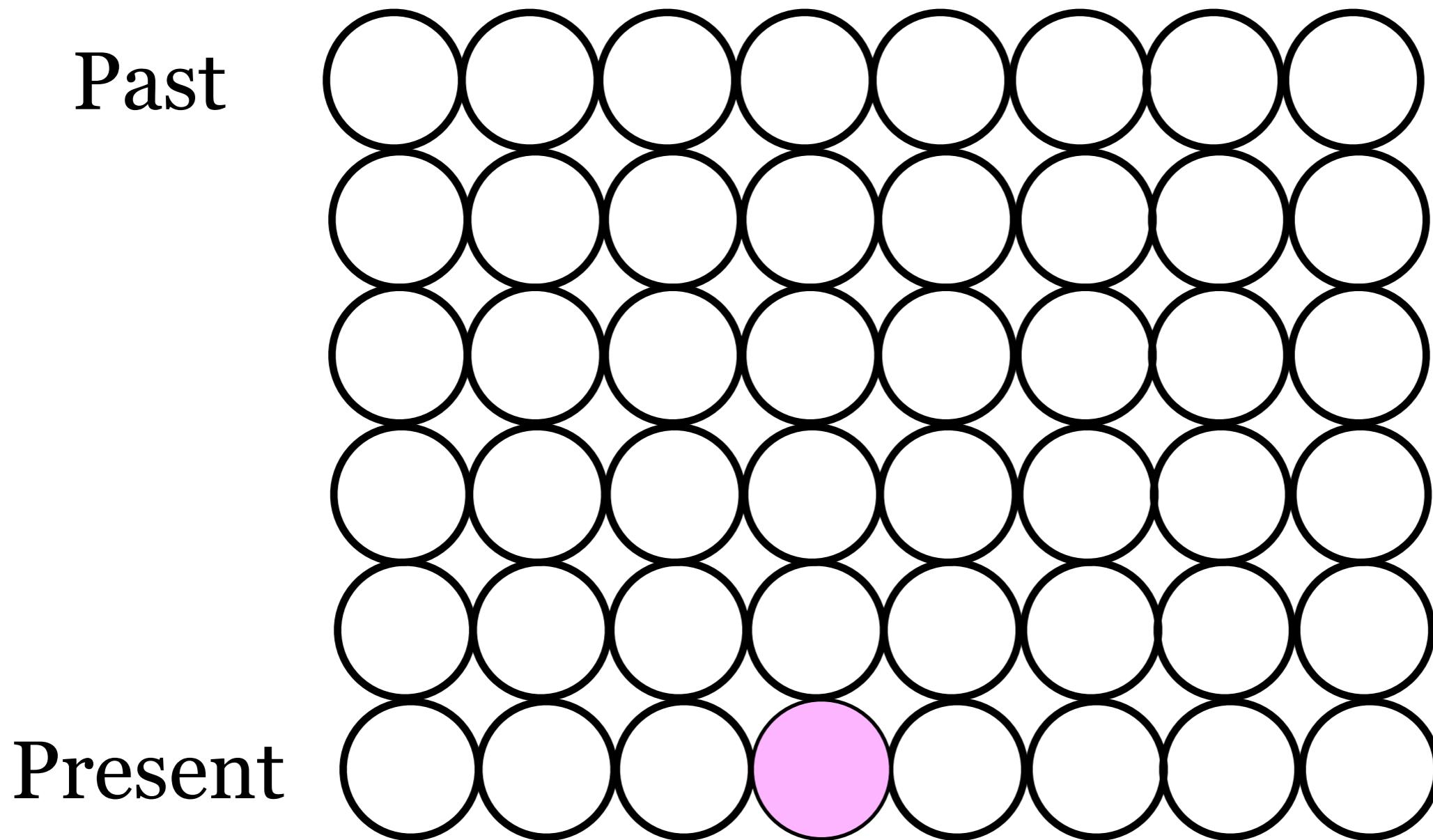
# Coalescent model within I

## population



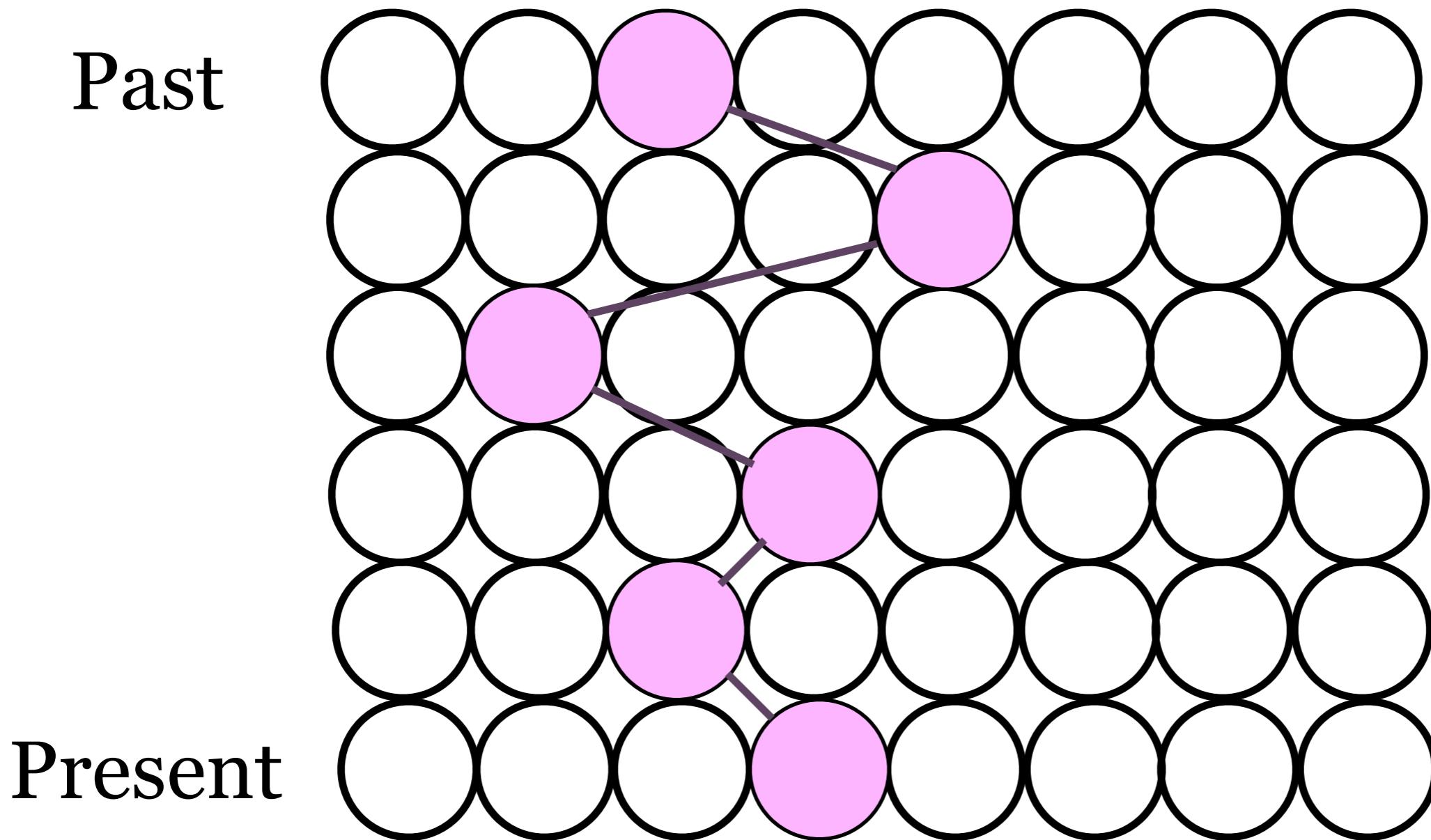
**How many generations do we have to wait for these two individuals to reach a common ancestor?**

# Coalescent model within I population



**How many generations do we have to wait for these two individuals to reach a common ancestor?**

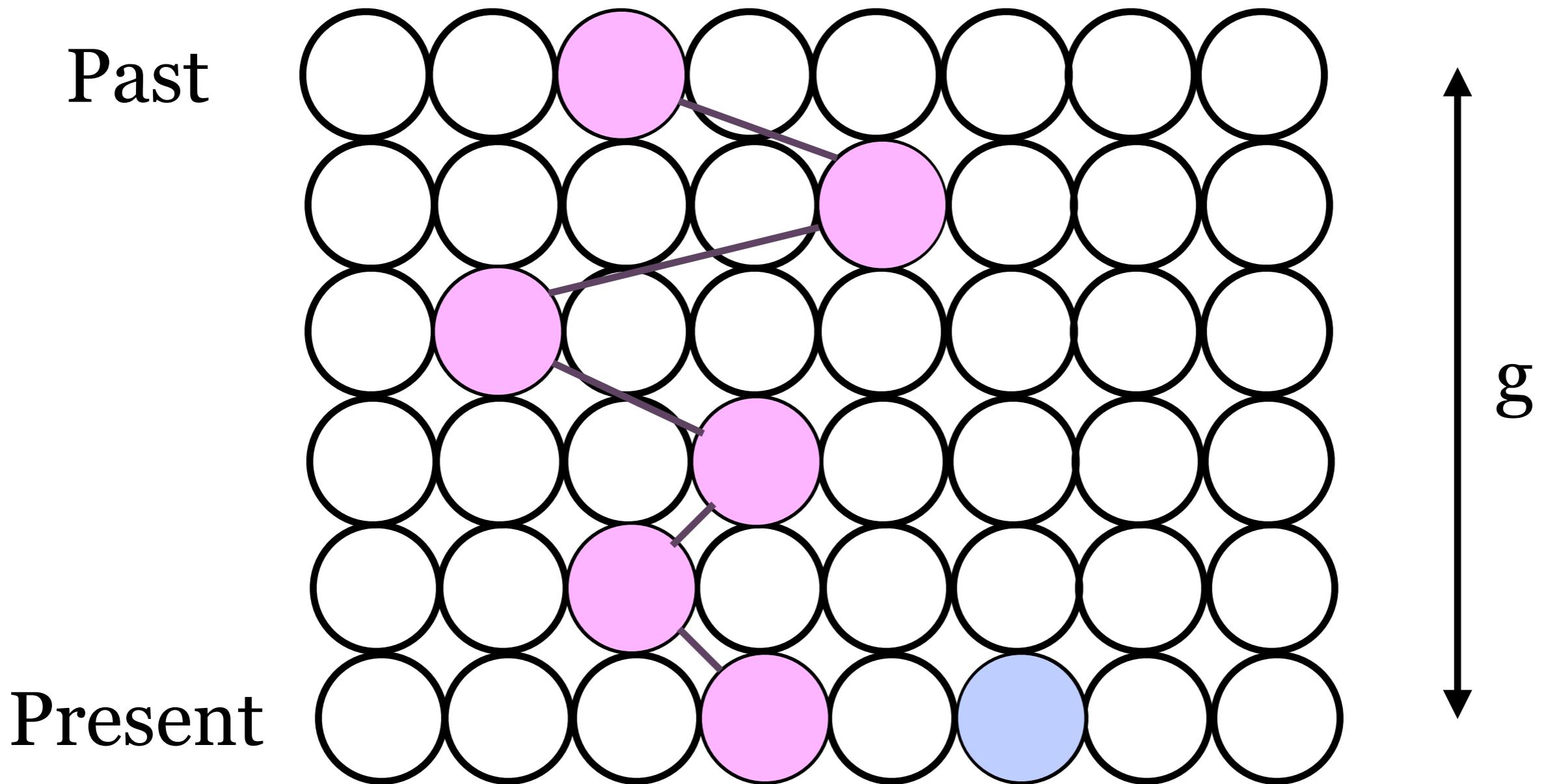
# Coalescent model within I population



**How many generations do we have to wait for these two individuals to reach a common ancestor?**

# Coalescent model within I

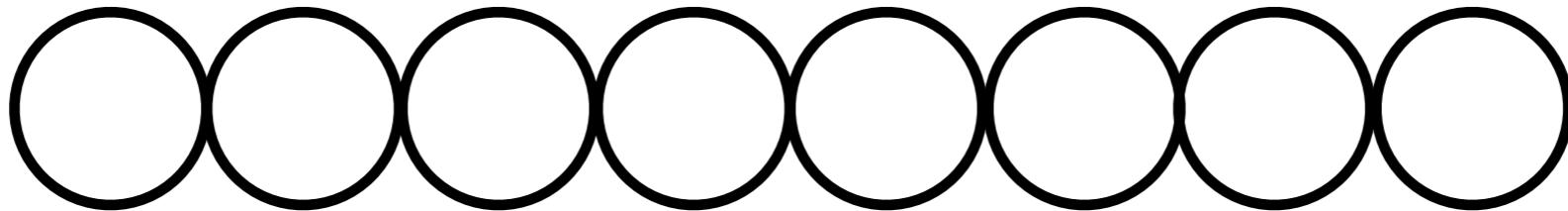
## population



**How many generations do we have to wait for these two individuals to reach a common ancestor?**

# Coalescent model within I population

Present



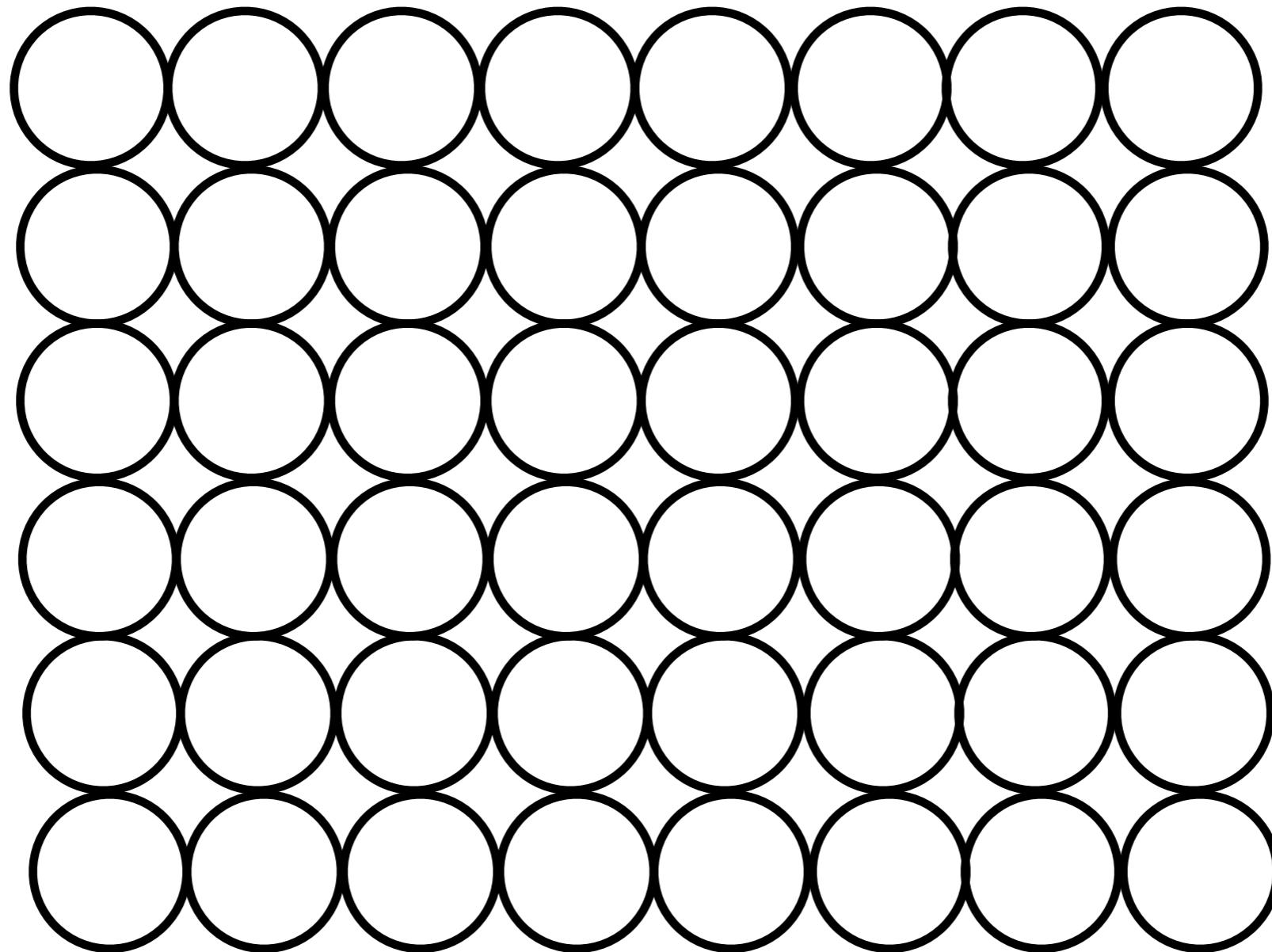
How many generations do we have  
to wait for these two individuals to  
reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

Past

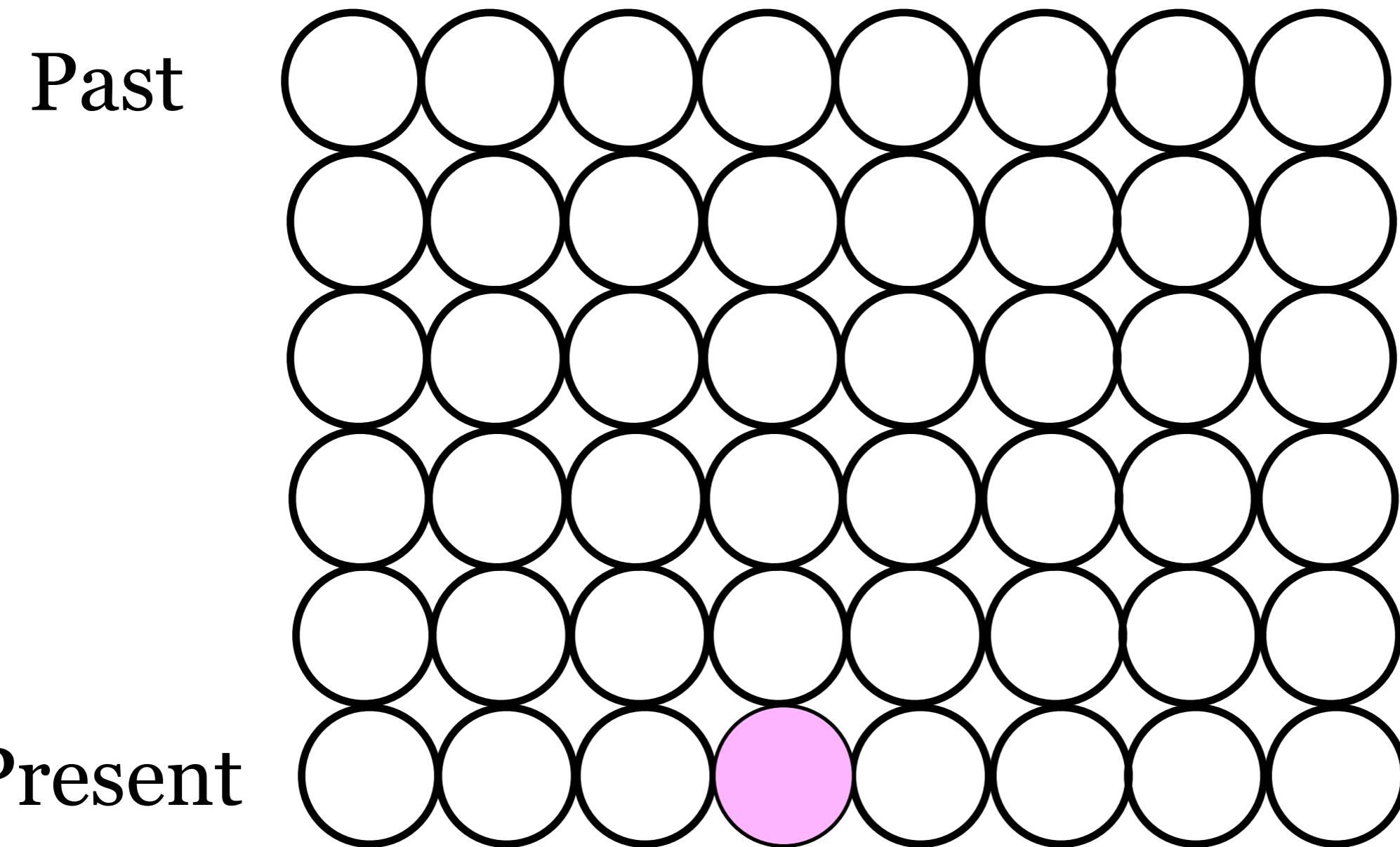


How many generations do we have  
to wait for these two individuals to  
reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

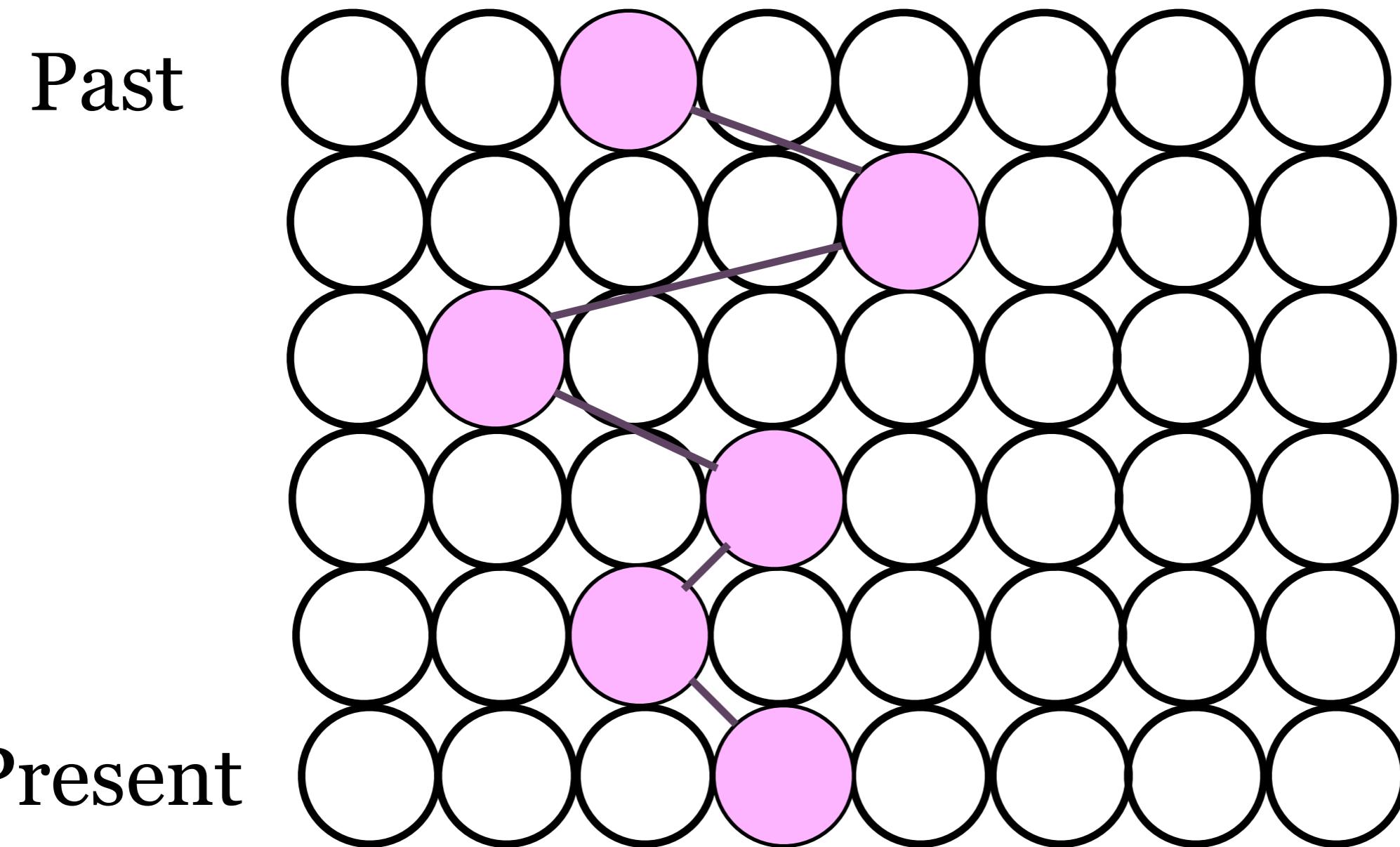


How many generations do we have to wait for these two individuals to reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

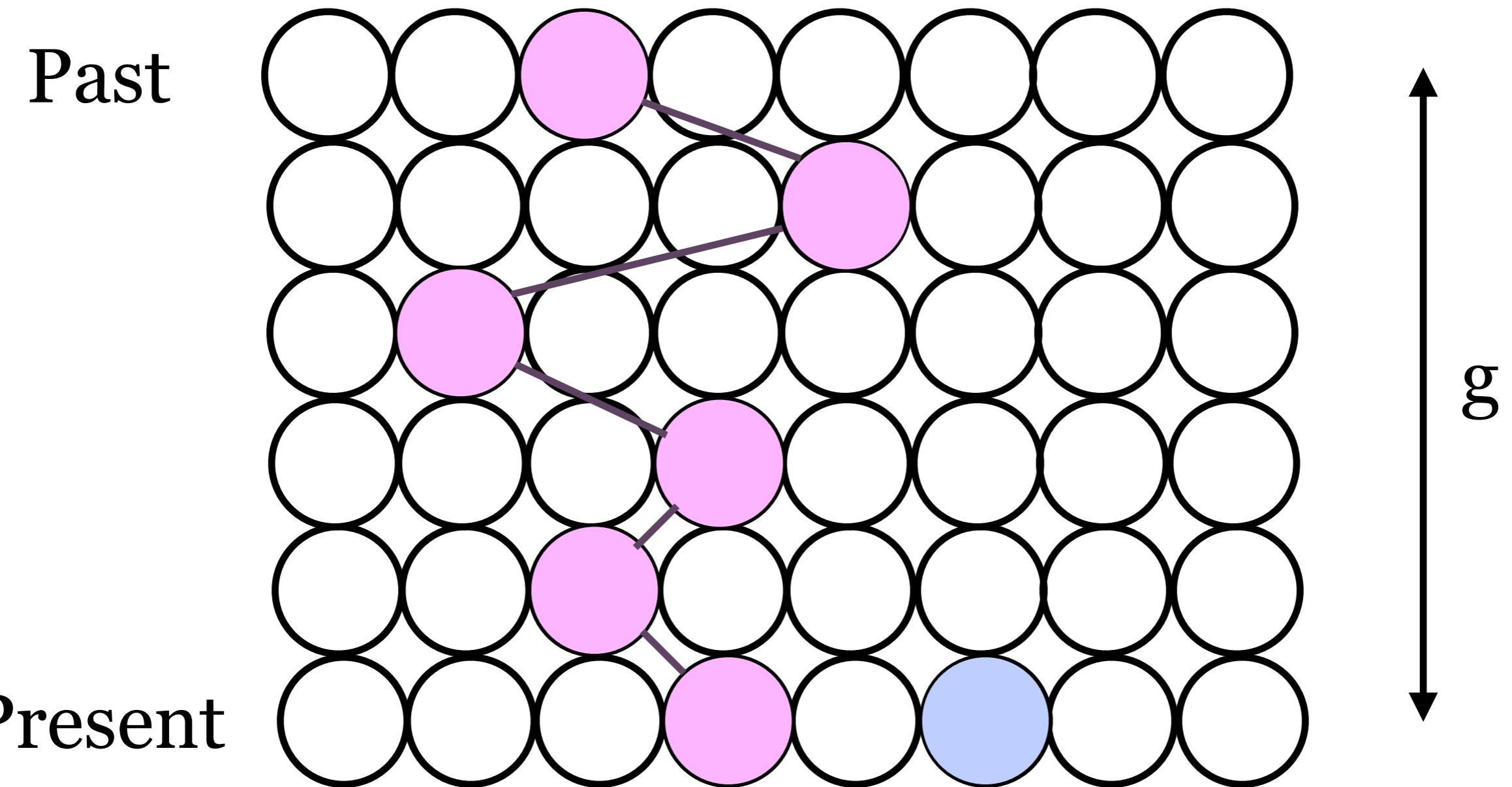


How many generations do we have to wait for these two individuals to reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population



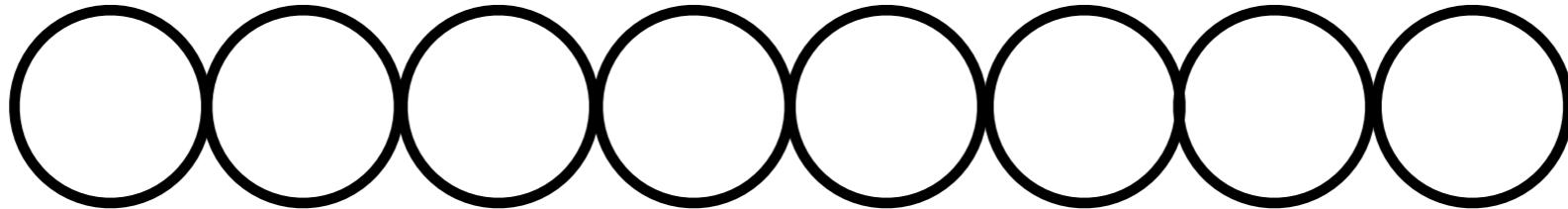
How many generations do we have to wait for these two individuals to reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

Present



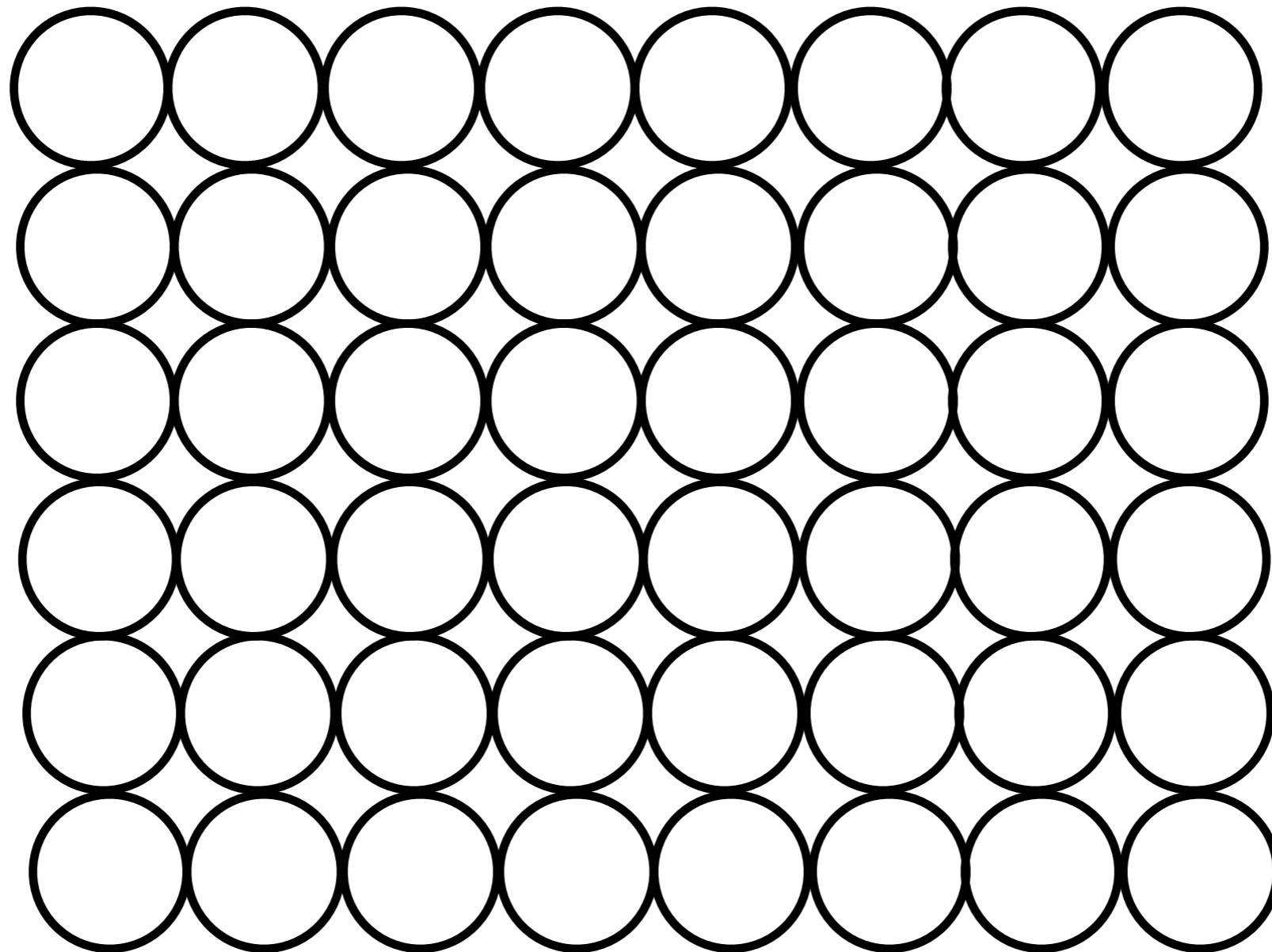
How many generations do we have  
to wait for these two individuals to  
reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

Past



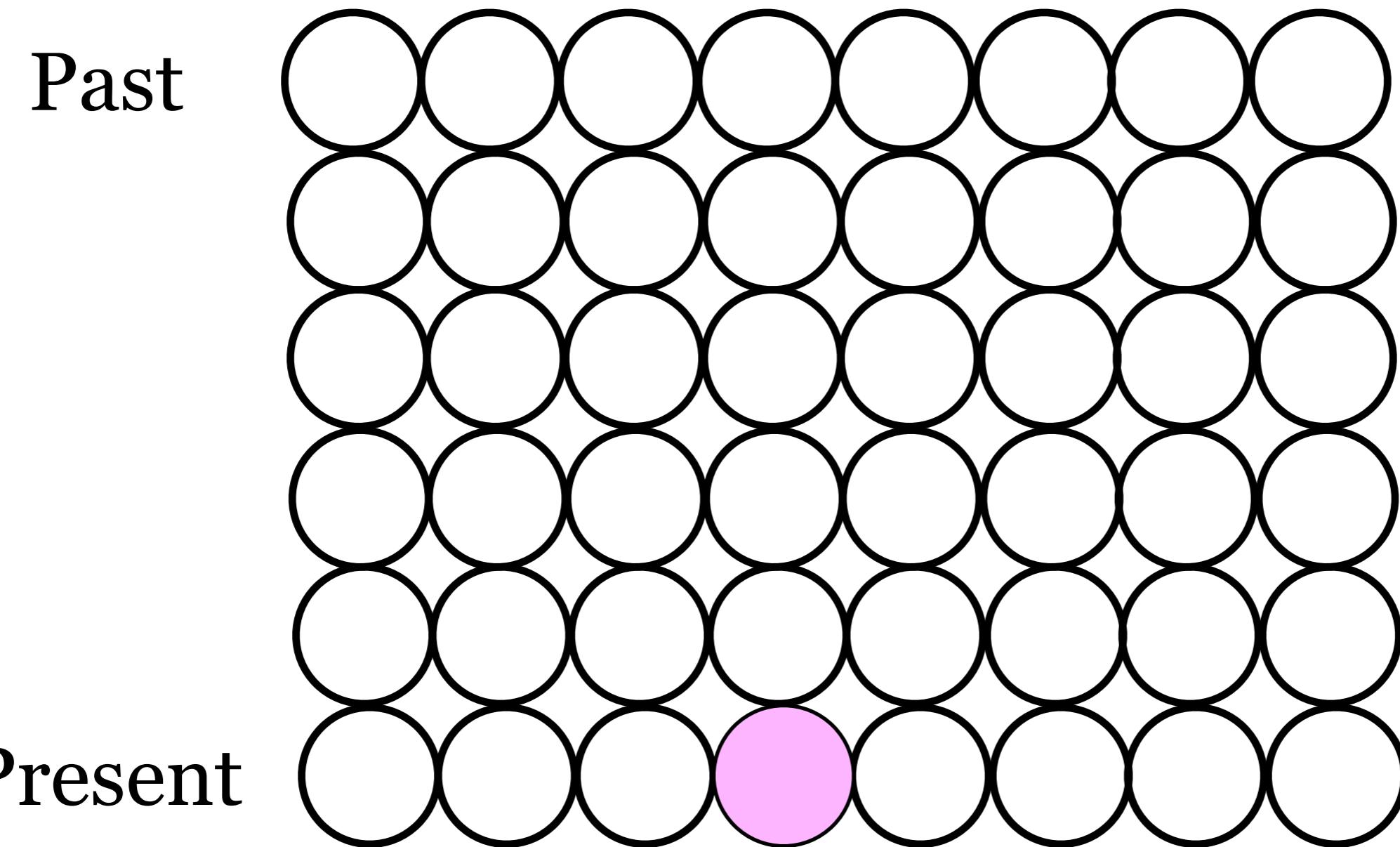
Present

How many generations do we have to wait for these two individuals to reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

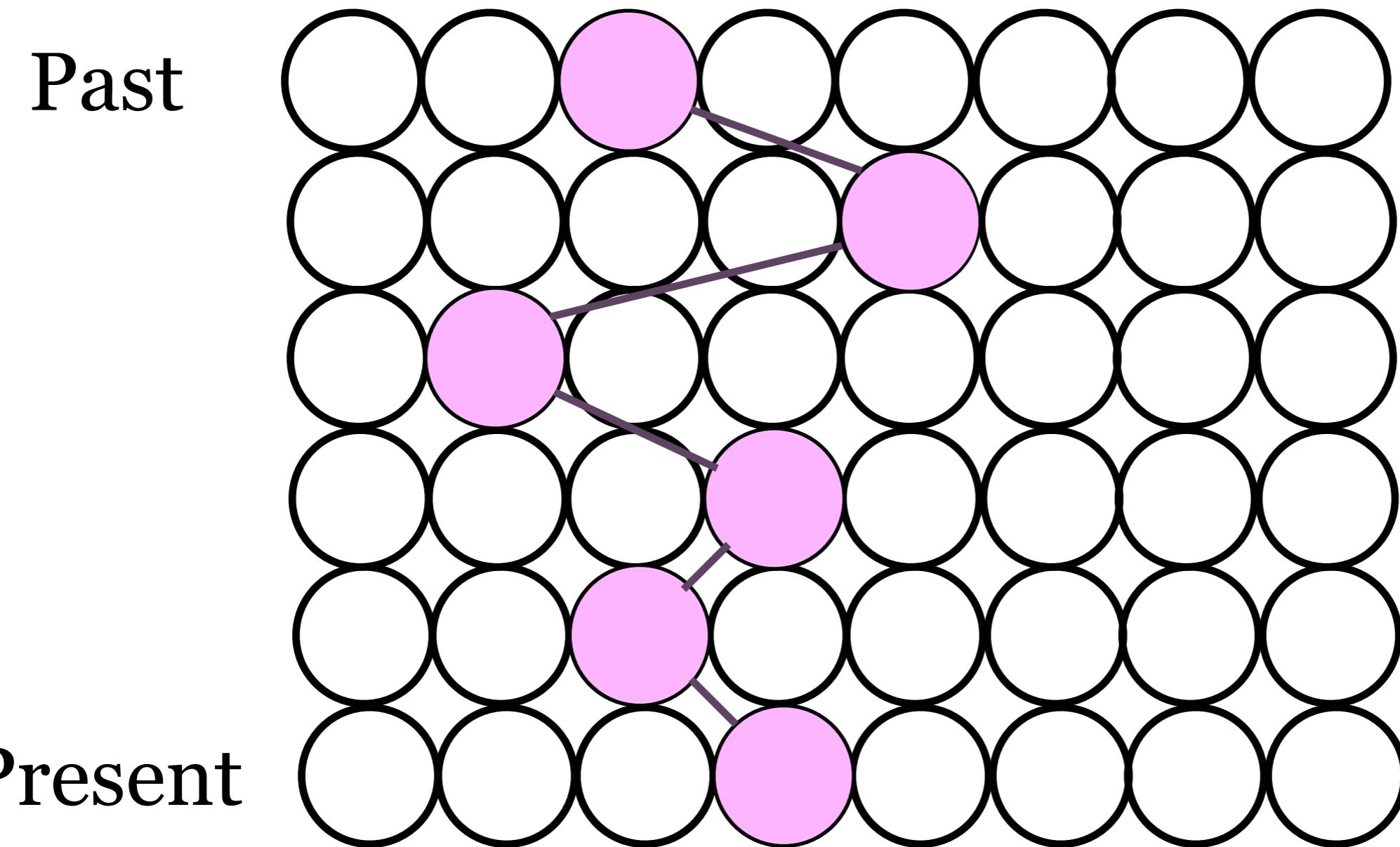


How many generations do we have to wait for these two individuals to reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

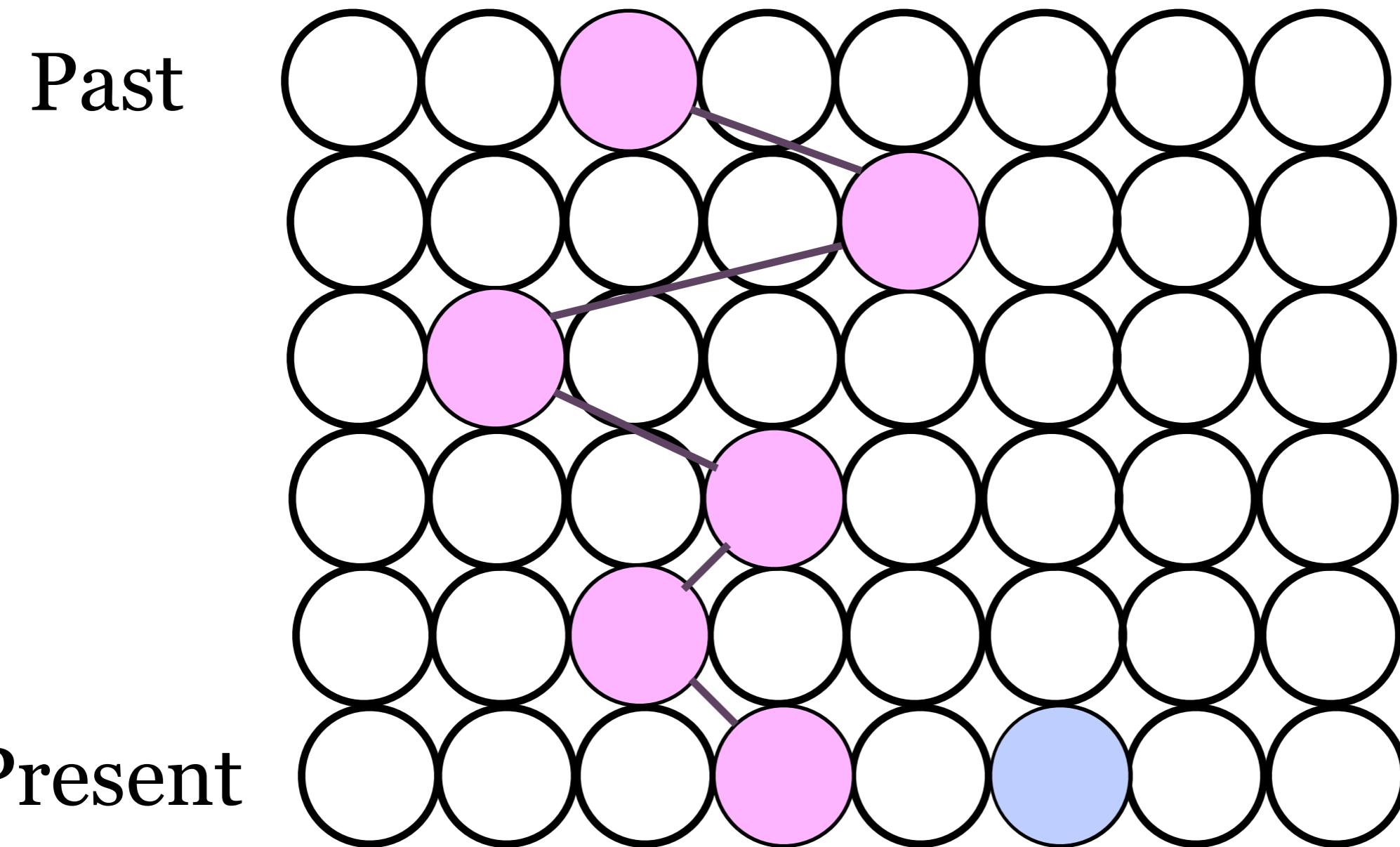


How many generations do we have  
to wait for these two individuals to  
reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population

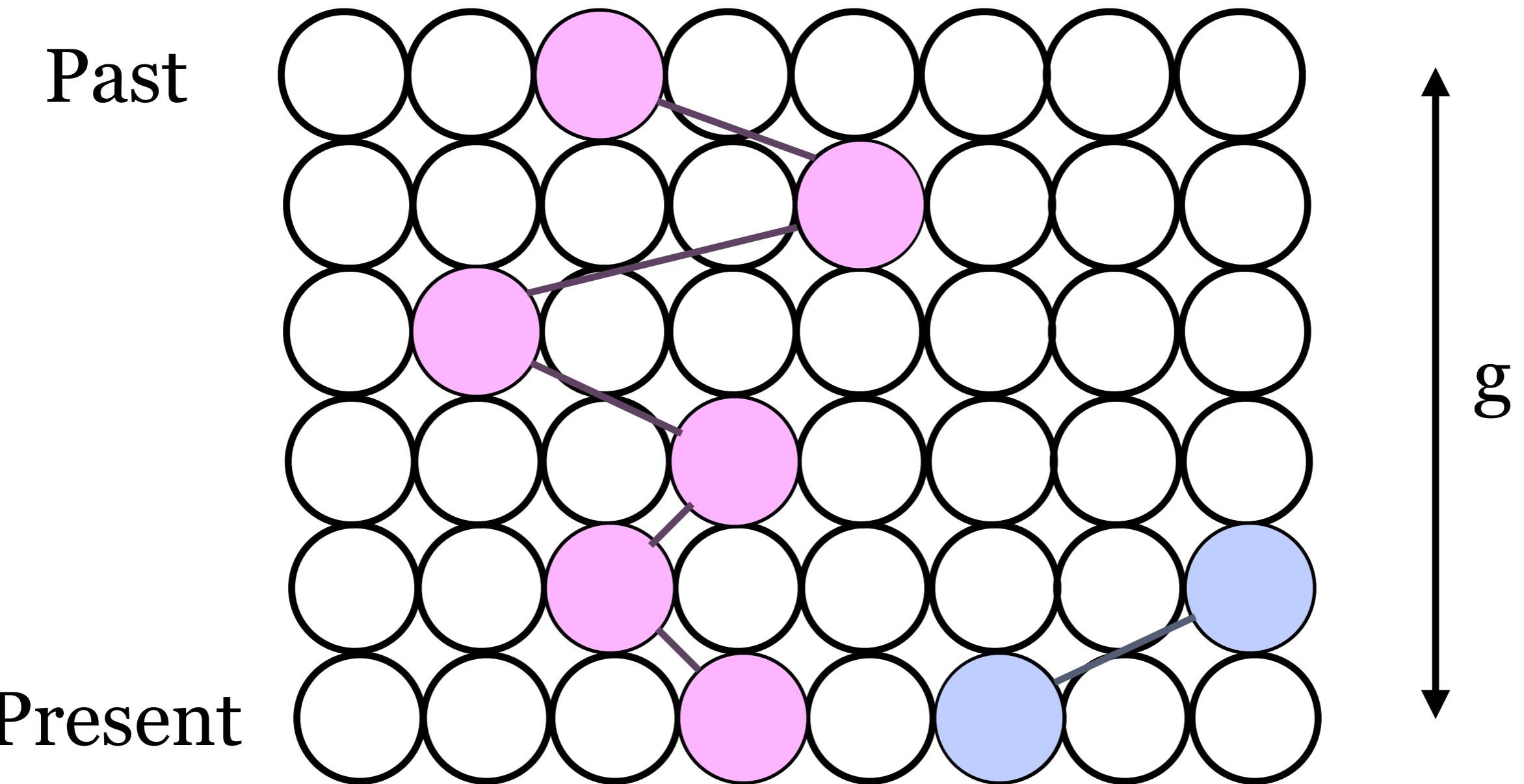


How many generations do we have to wait for these two individuals to reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

# Coalescent model within I population



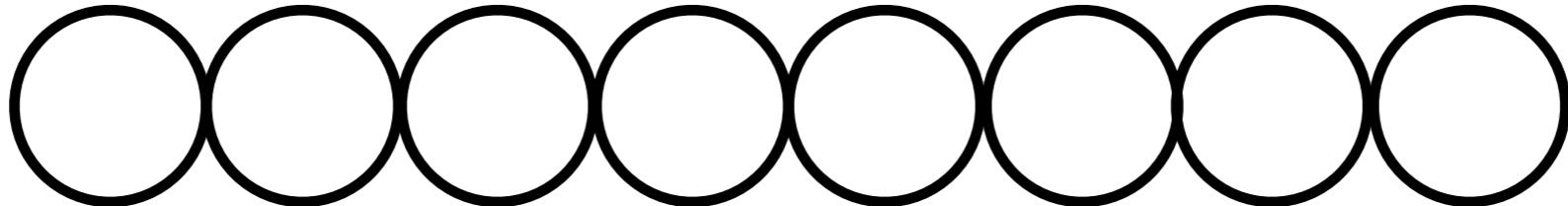
How many generations do we have  
to wait for these two individuals to  
reach a common ancestor?

$$P(\text{coalesce}) = \frac{1}{N}$$

$$P(\text{no coalesce}) = 1 - \frac{1}{N}$$

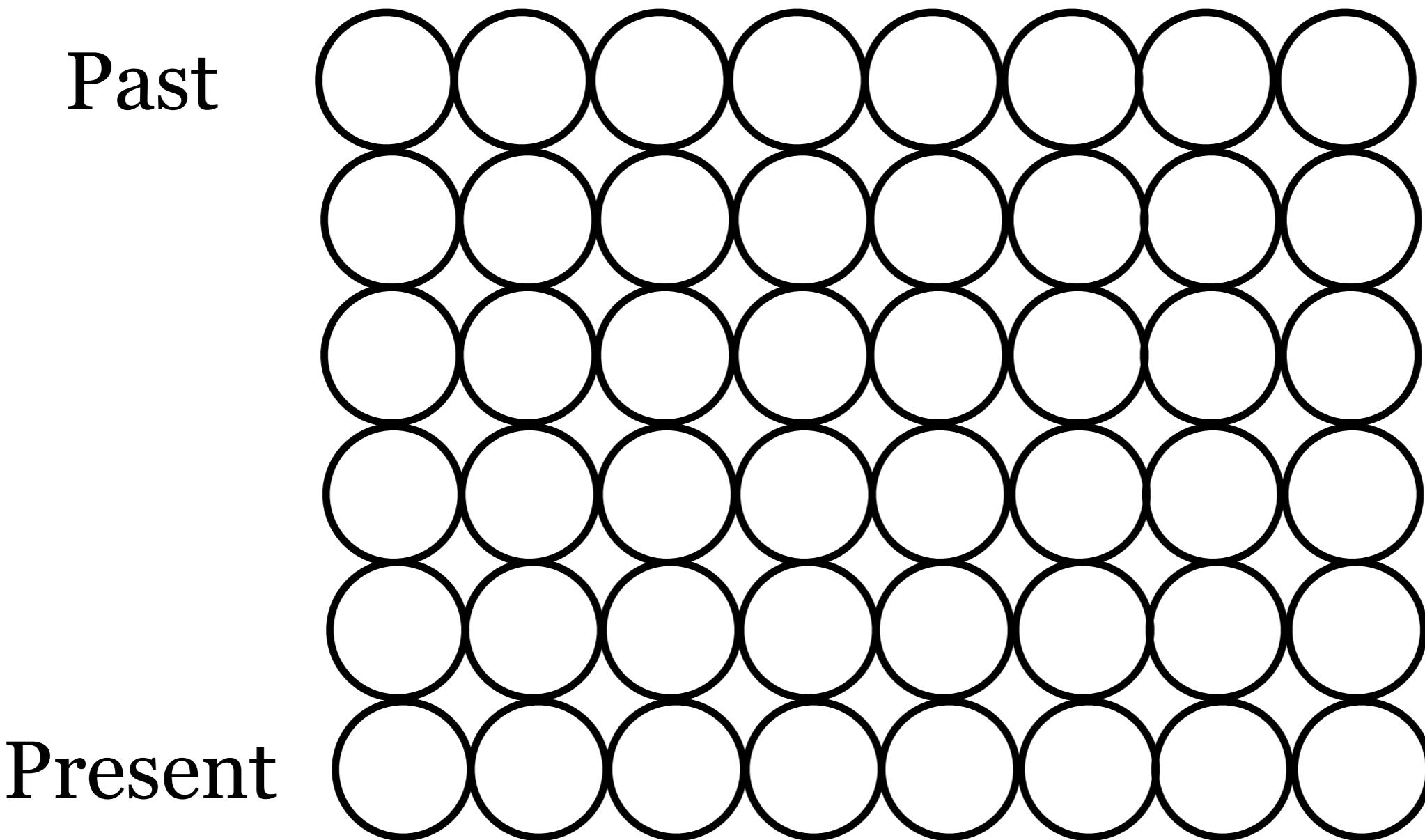
# Coalescent model within I population

Present



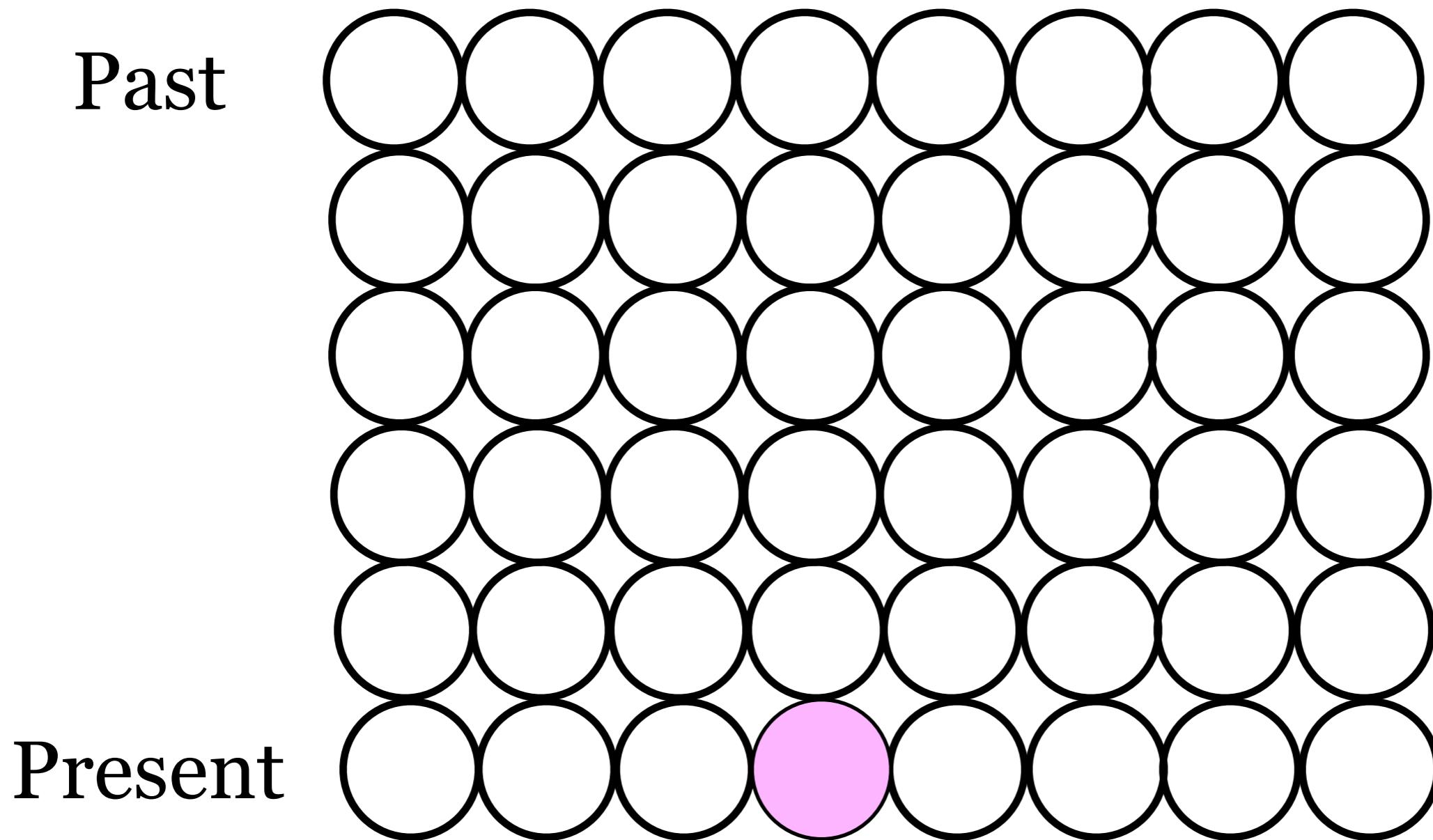
# Coalescent model within I

population



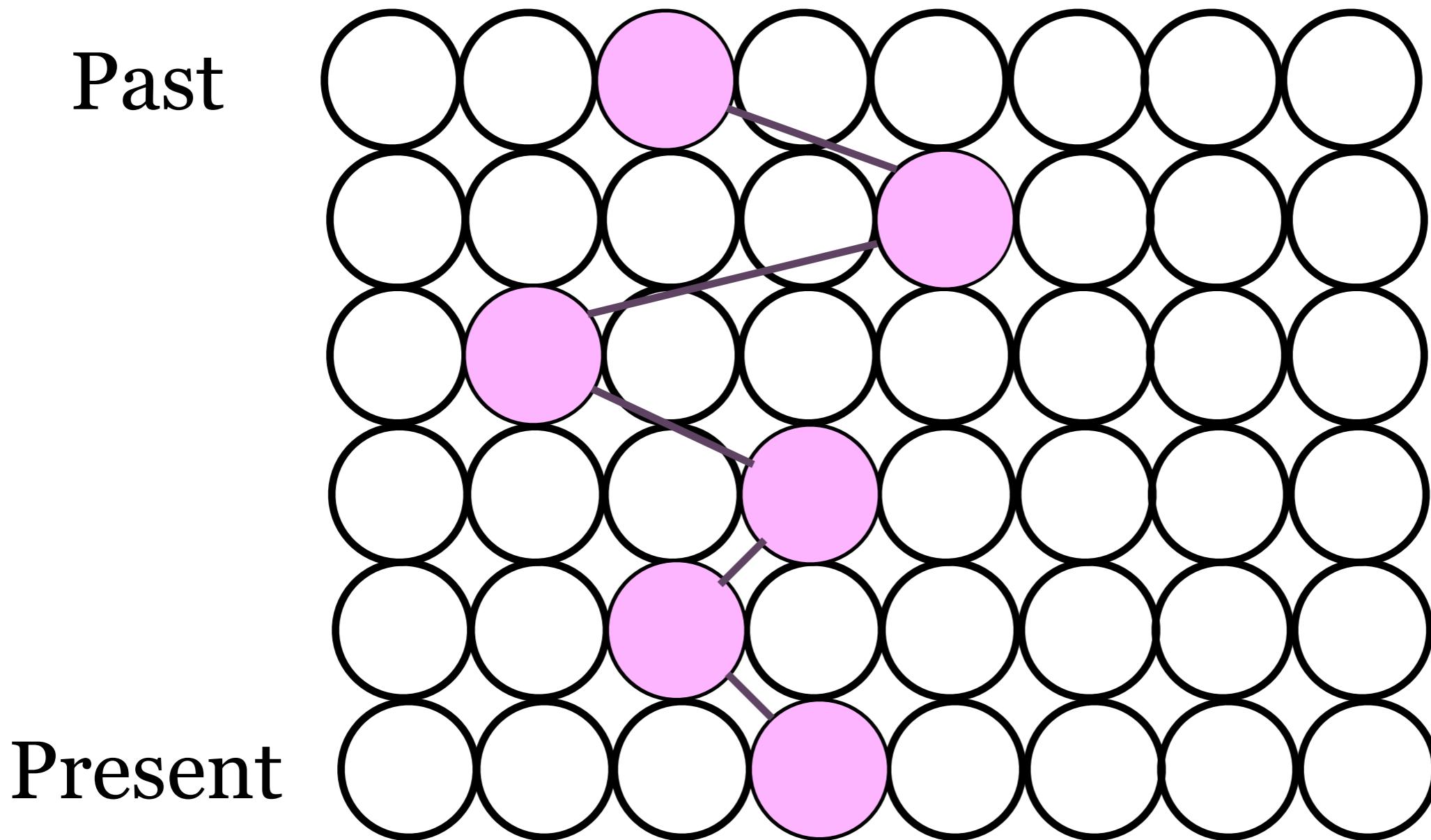
# Coalescent model within I

population



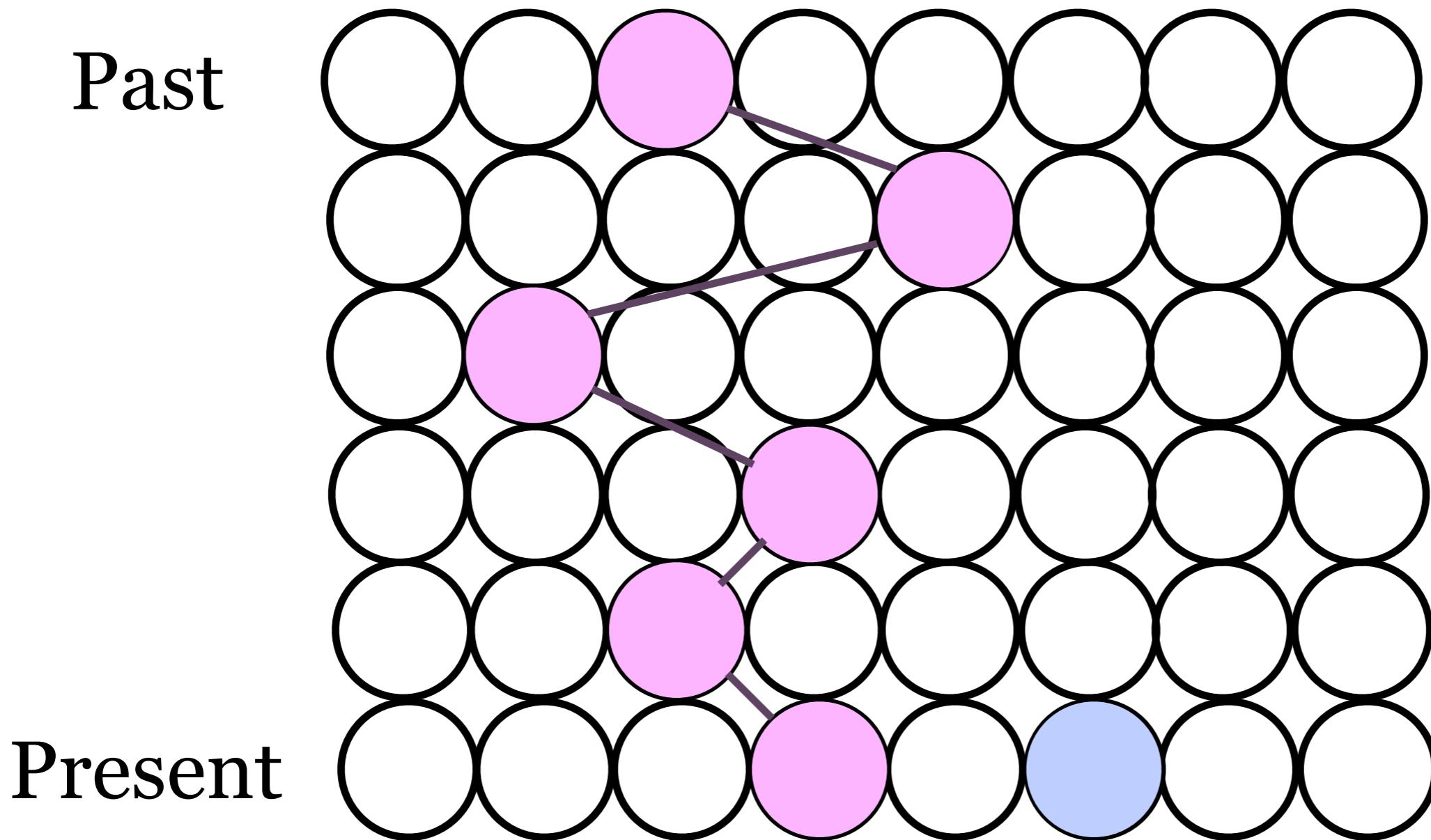
# Coalescent model within I

population



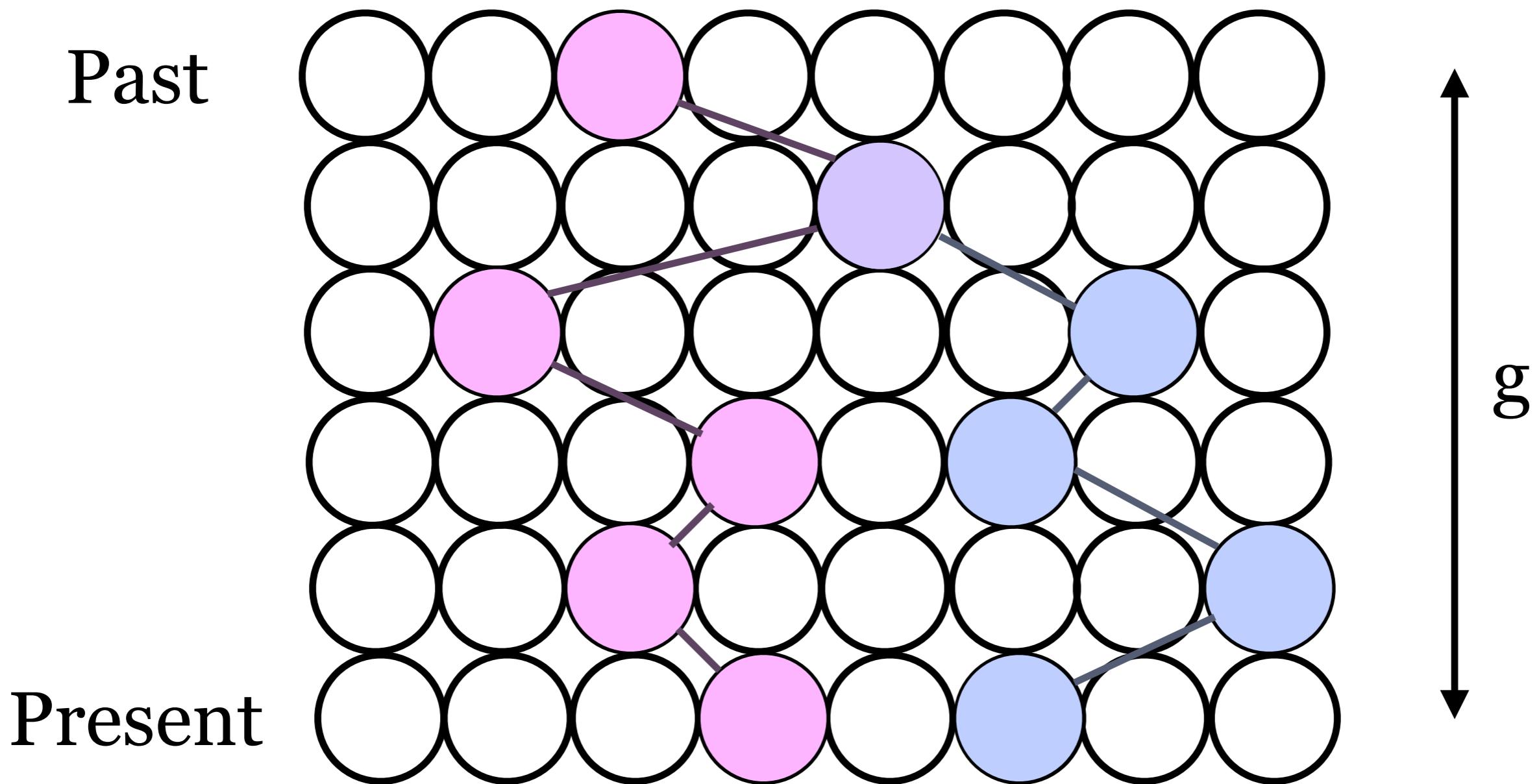
# Coalescent model within I

population



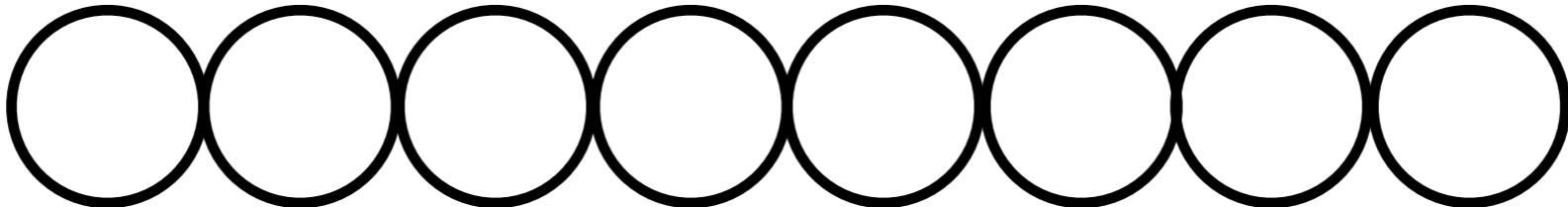
# Coalescent model within I

population



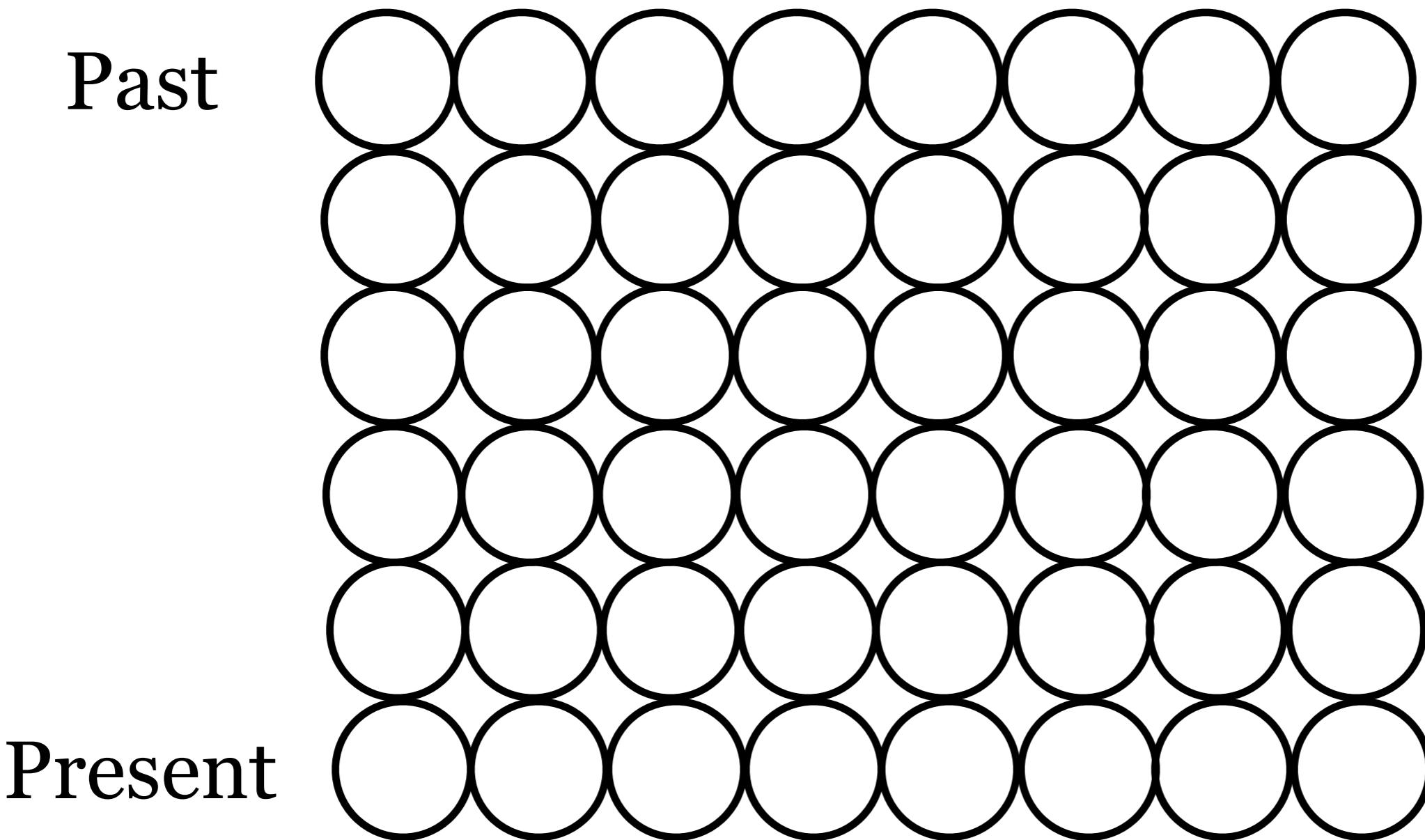
# Coalescent model within I population

Present



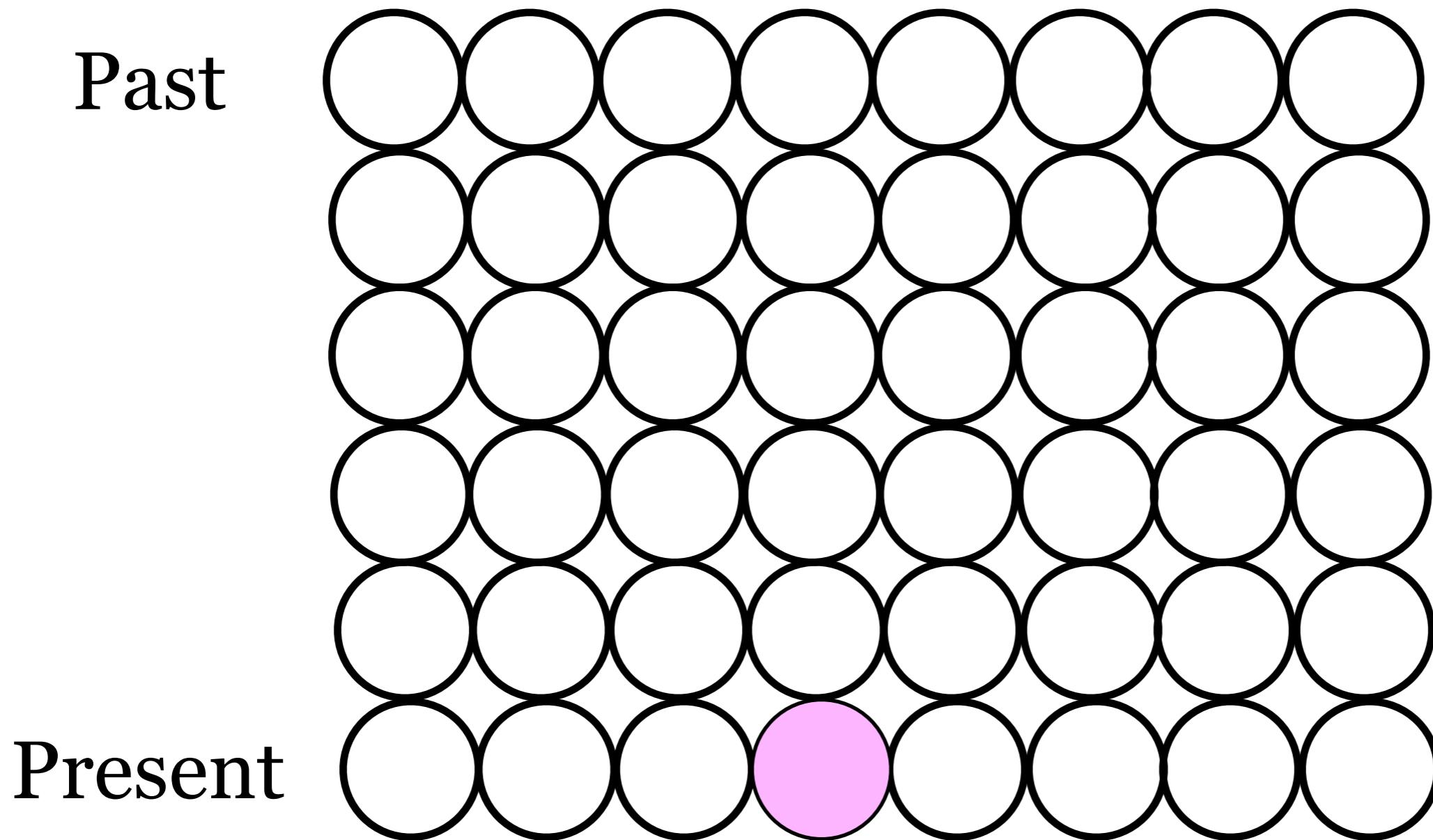
# Coalescent model within I

population



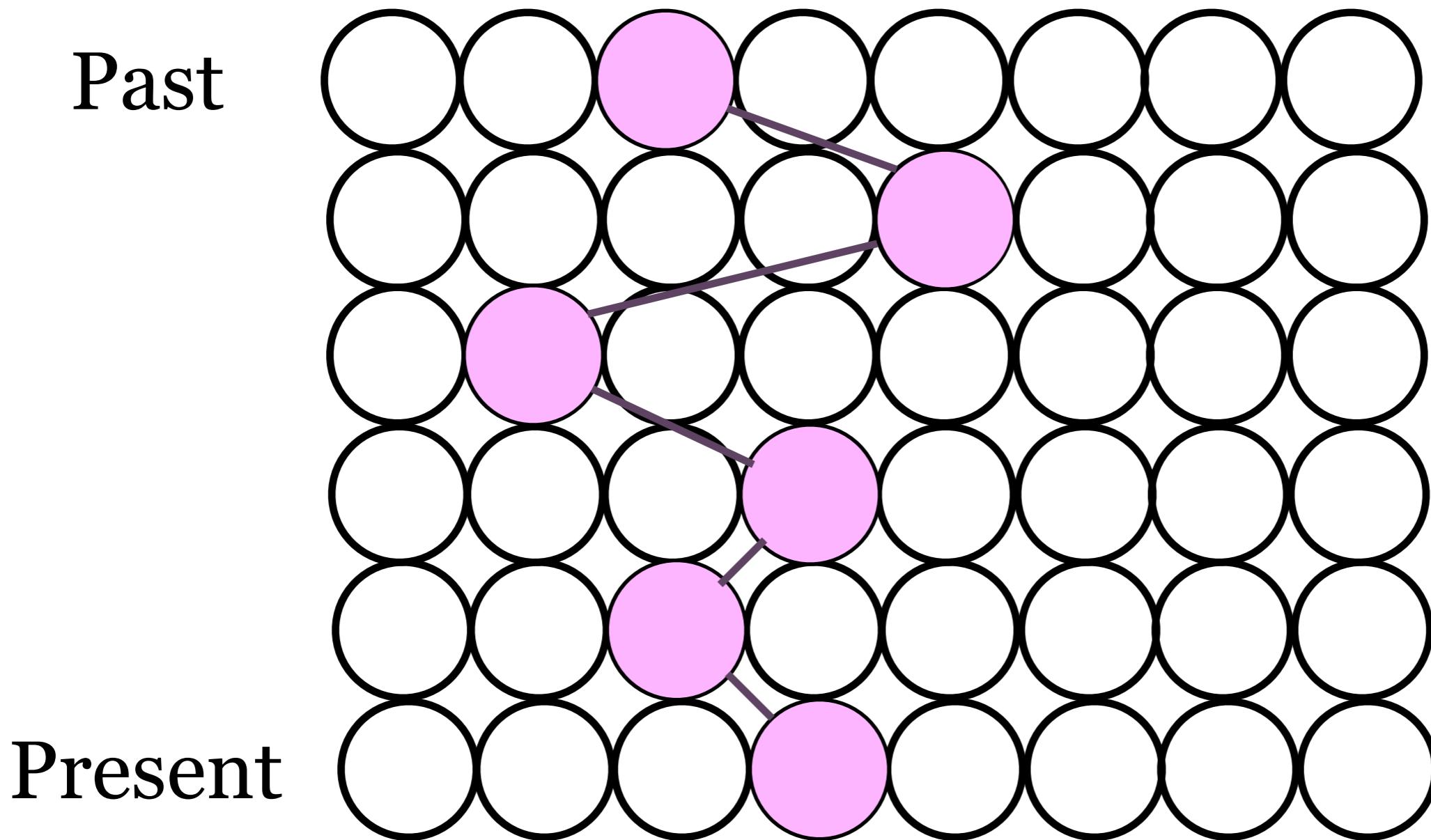
# Coalescent model within I

population



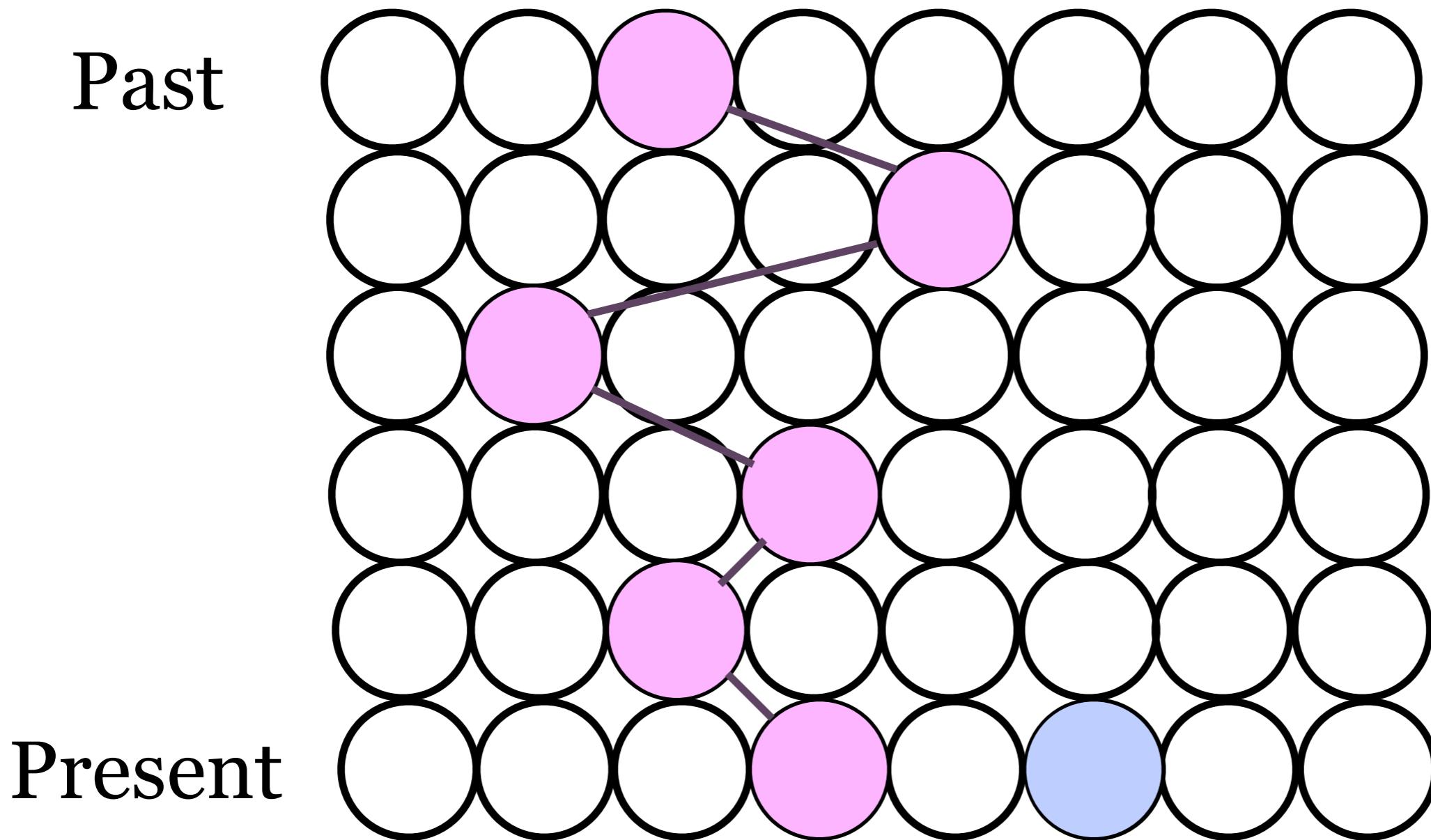
# Coalescent model within I

population



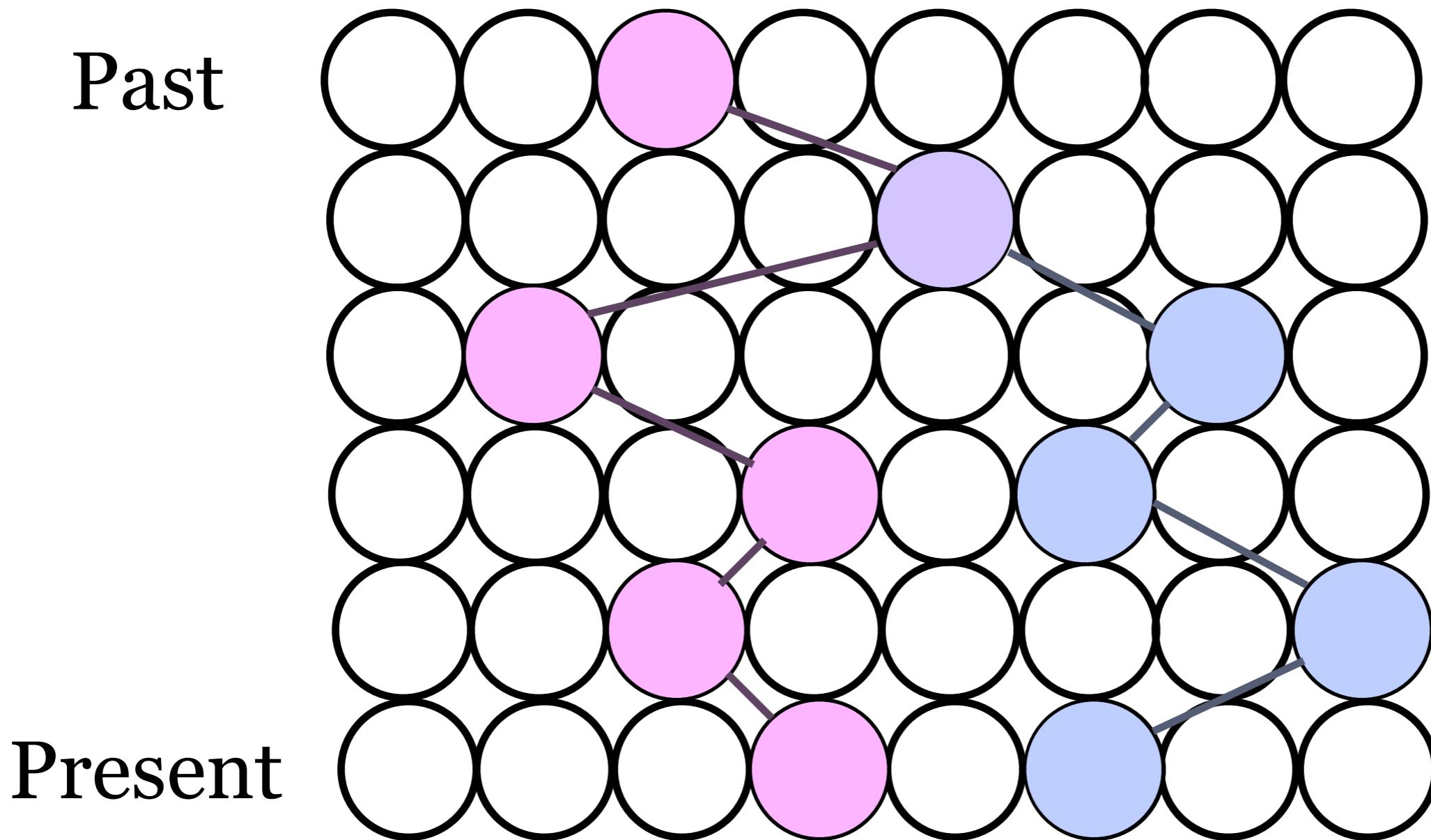
# Coalescent model within I

population



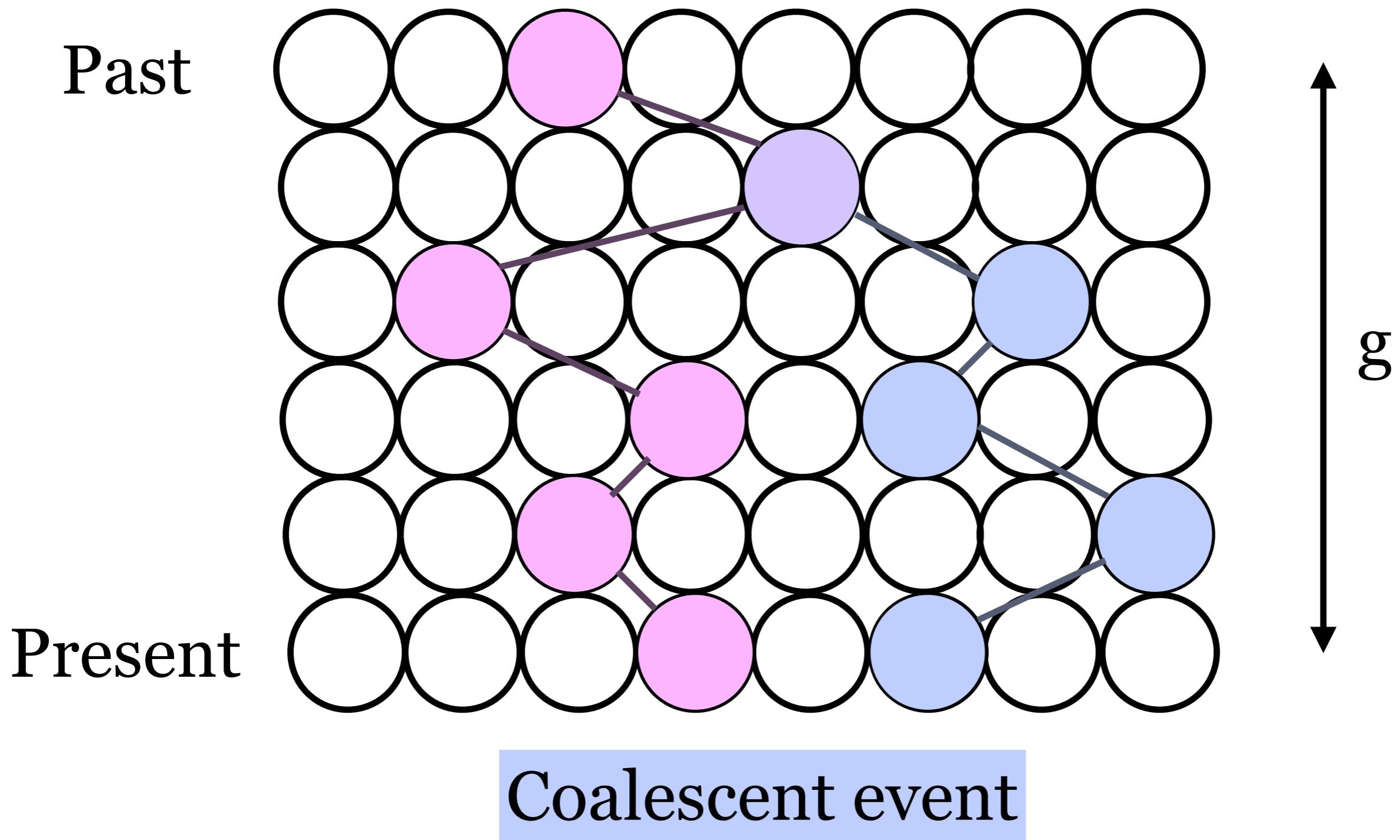
# Coalescent model within I

population

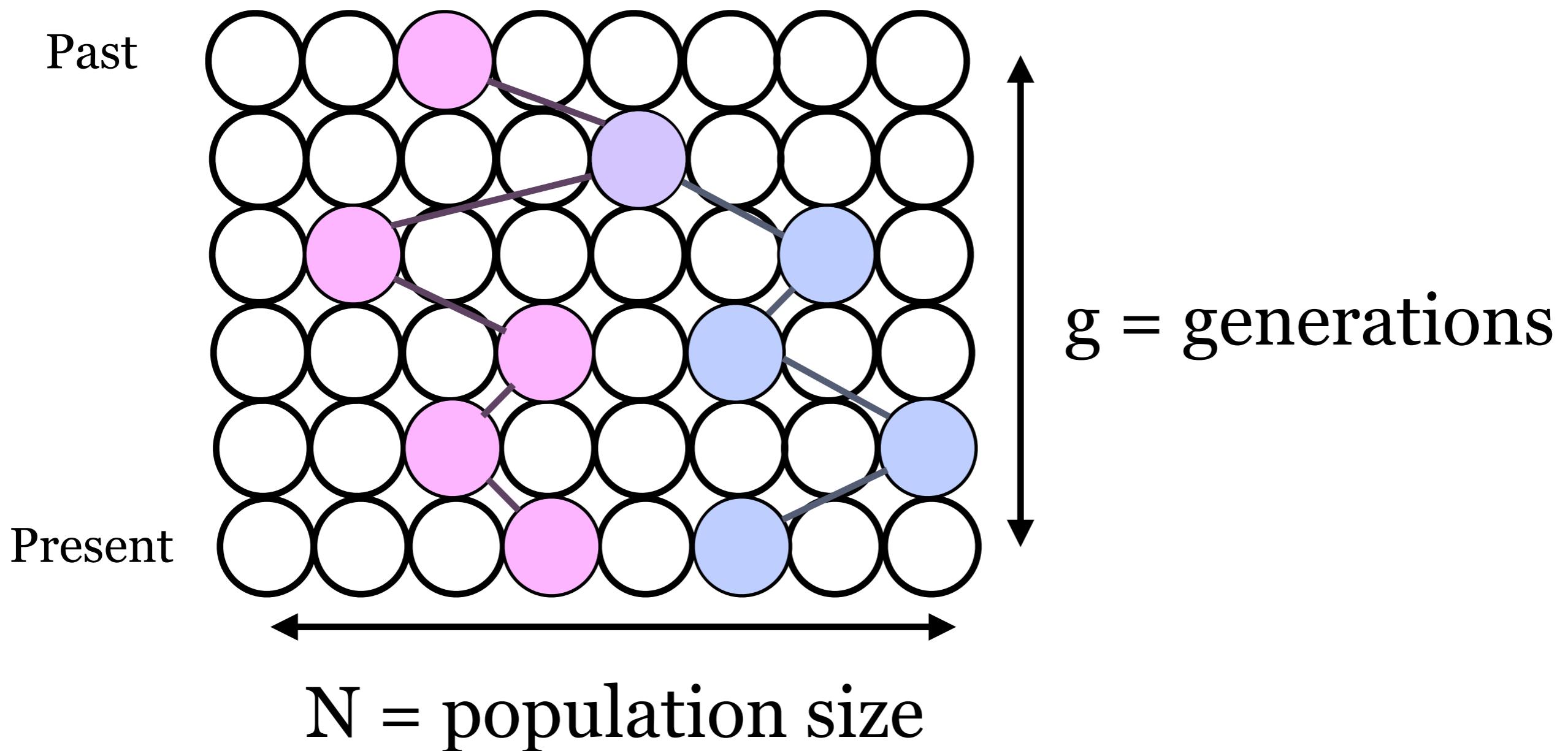


# Coalescent model within I

population



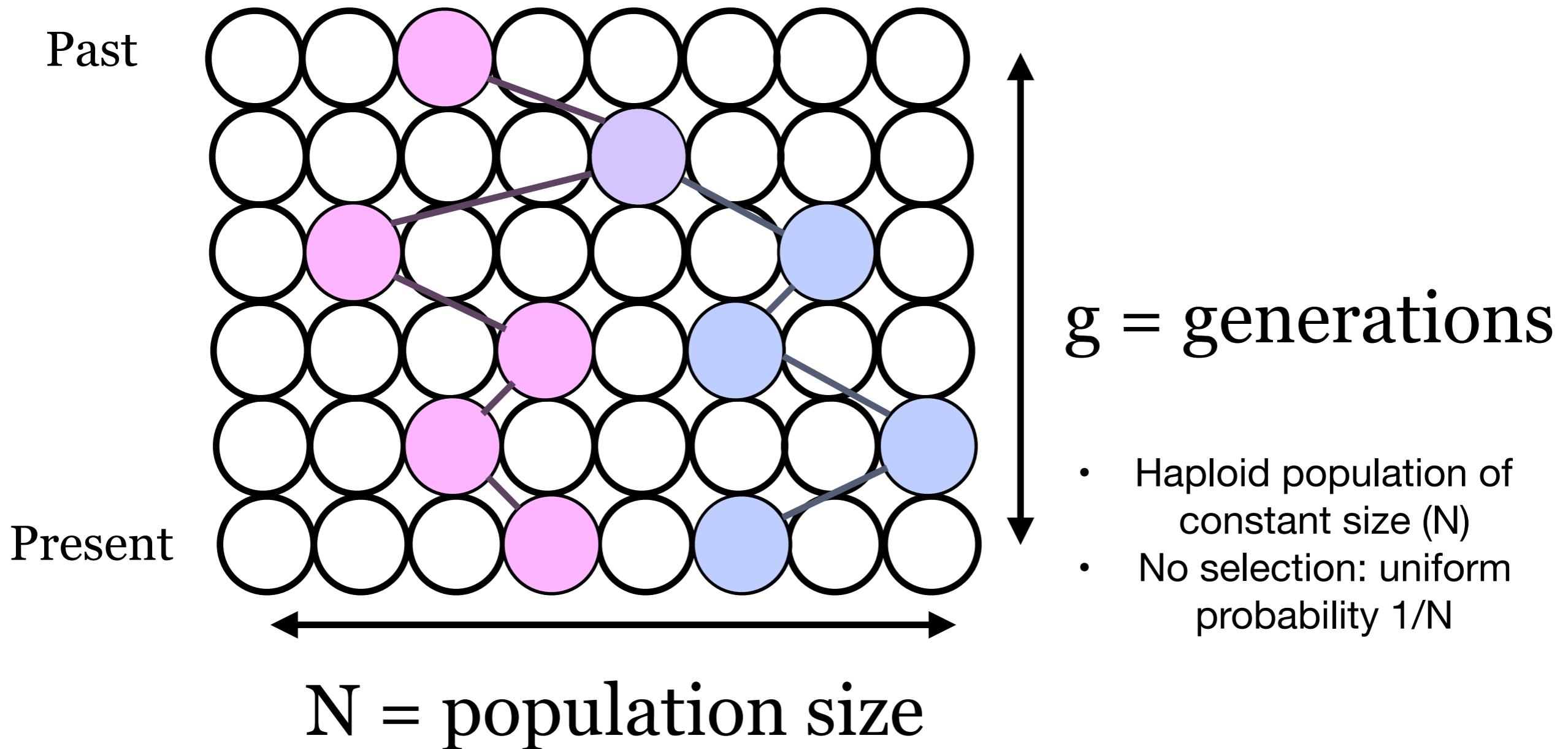
# Coalescent model within 1 population



Probability of no coalescence in  $g$  generations:

$$\left(1 - \frac{1}{N}\right)^g$$
$$t = g/N \Rightarrow \left(1 - \frac{t}{Nt}\right)^{Nt} \xrightarrow[N \rightarrow \infty]{} e^{-t}$$

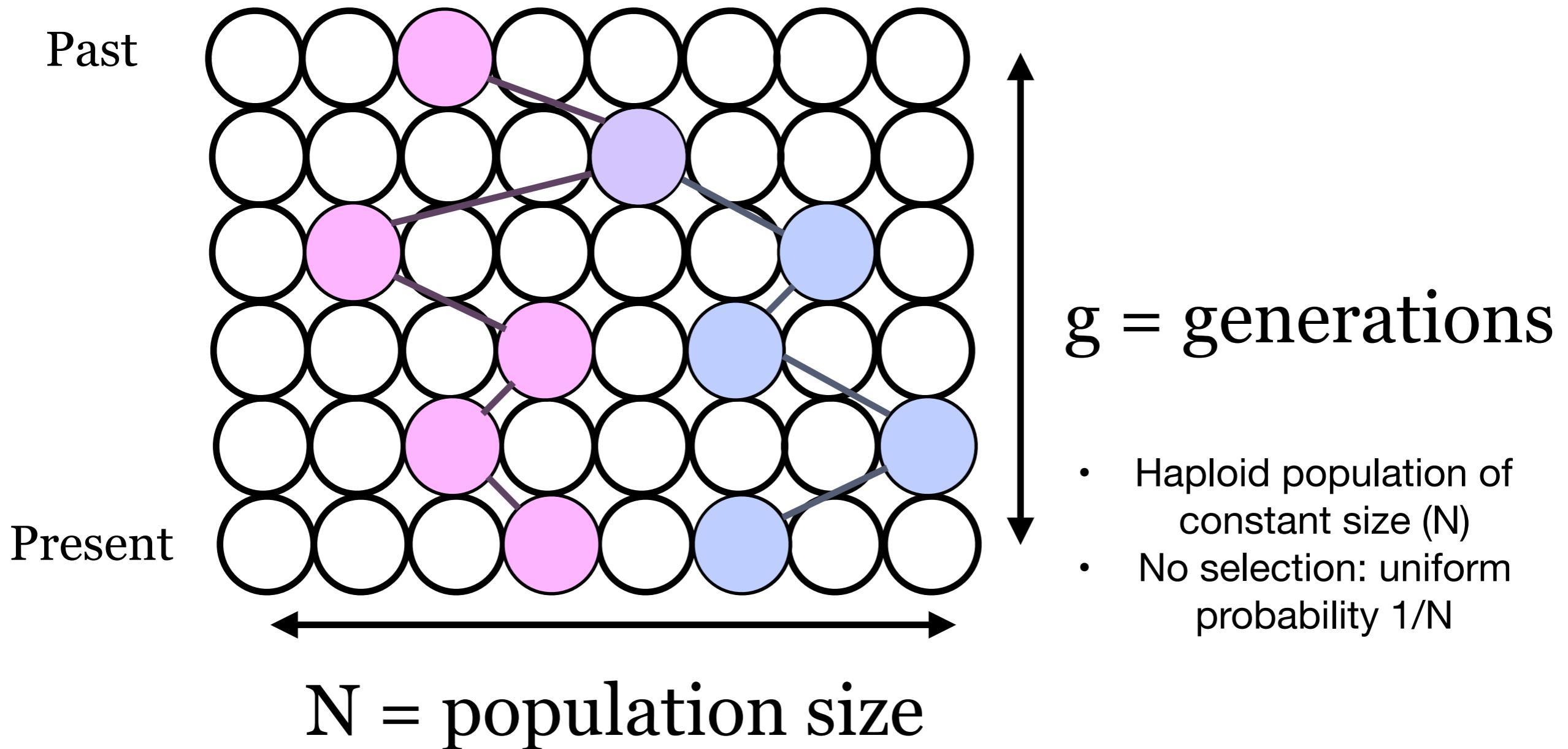
# Coalescent model within 1 population



Probability of no coalescence in  $g$  generations:

$$\left(1 - \frac{1}{N}\right)^g$$
$$t = g/N \Rightarrow \left(1 - \frac{t}{Nt}\right)^{Nt} \xrightarrow[N \rightarrow \infty]{} e^{-t}$$

# Coalescent model within 1 population

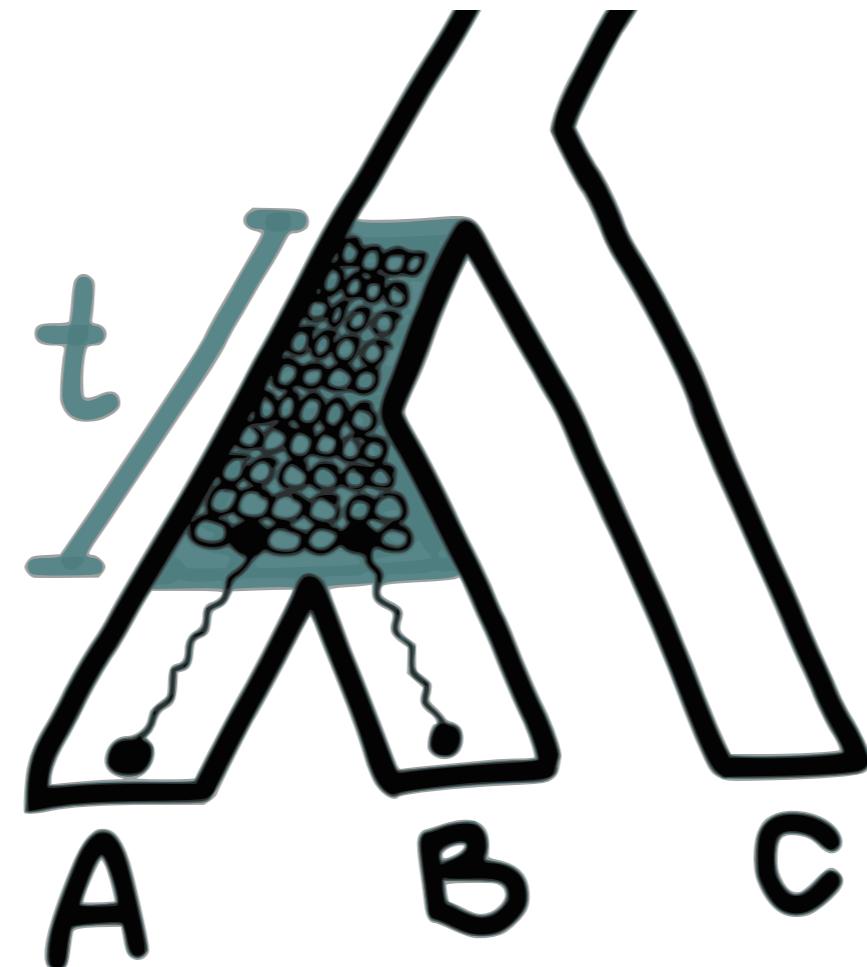


Probability of no coalescence in  $g$  generations:

$$t = g/N \Rightarrow \left(1 - \frac{t}{Nt}\right)^{Nt} \xrightarrow[N \rightarrow \infty]{} e^{-t}$$
$$\left(1 - \frac{1}{N}\right)^g$$

$\boxed{2N}$

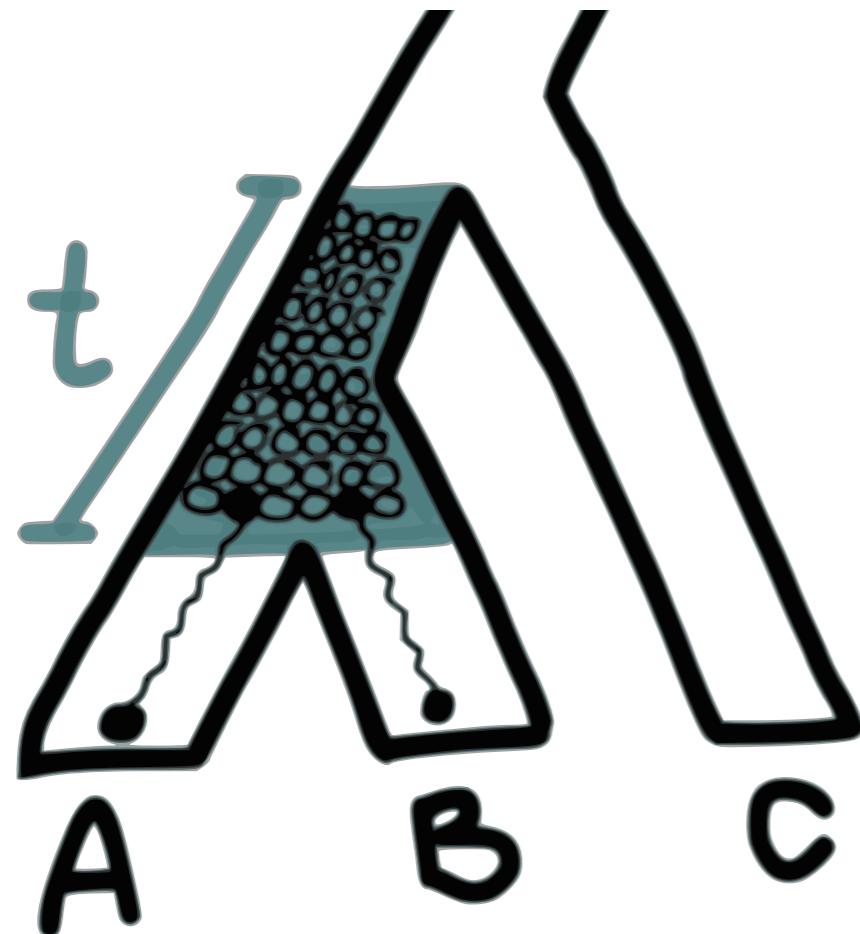
# Multispecies coalescent on a tree



$$P(T > t) = e^{-t}$$

$$T = \frac{g}{N} \text{ coalescent units} \sim \text{Exp}(1)$$

# Multispecies coalescent on a tree

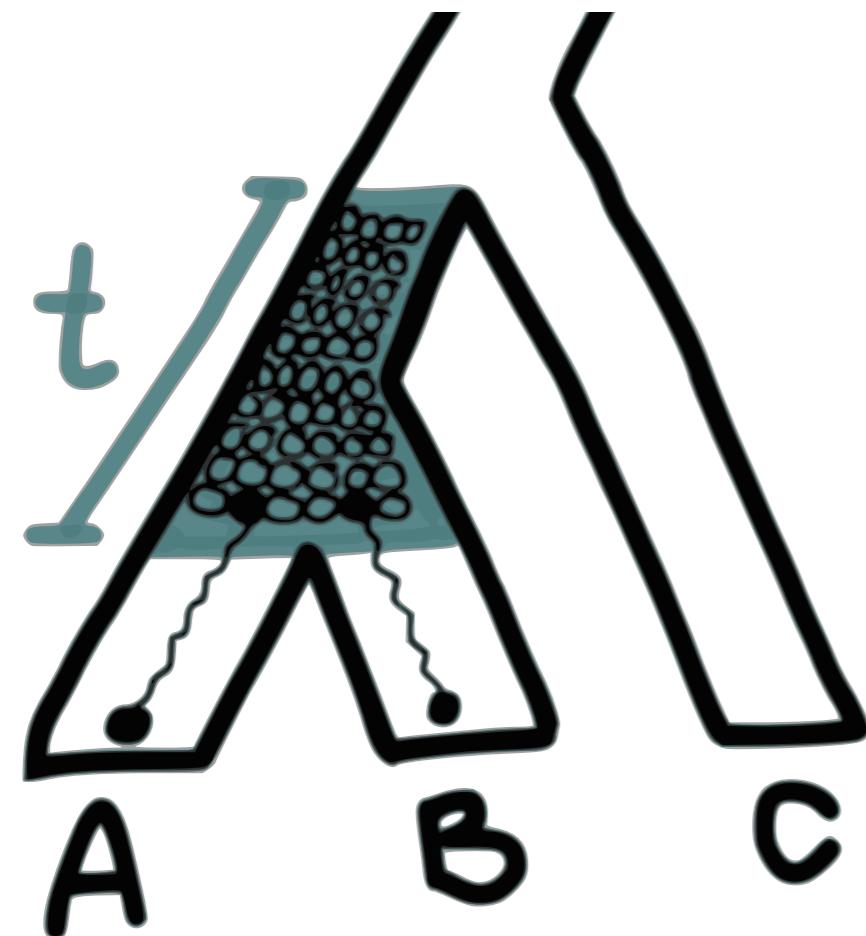


$$P(\text{ } \text{ } \text{ } \text{ } \text{ } ) =$$

A probability expression  $P(\text{ } \text{ } \text{ } \text{ } \text{ } ) =$  followed by a phylogenetic tree with three tips labeled A, B, and C.

$$P(T > t) = e^{-t}$$

# Multispecies coalescent on a tree

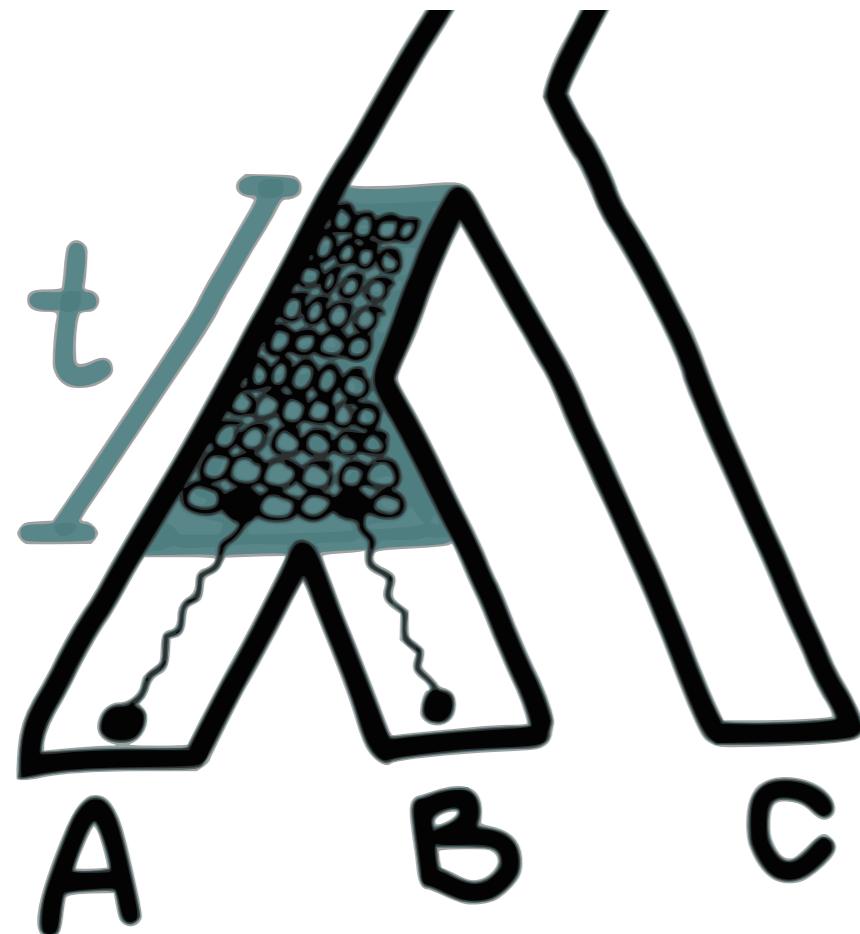


$$P(\text{ } \text{ } \text{ } \text{ } \text{ } \text{ } ) = 1 - e^{-t}$$

The probability of finding a specific tree topology (in this case, a star-like tree where all species A, B, and C share a single common ancestor) is given by the formula  $P(\text{ } \text{ } \text{ } \text{ } \text{ } \text{ } ) = 1 - e^{-t}$ , where  $t$  is the time since the most recent common ancestor.

$$P(T > t) = e^{-t}$$

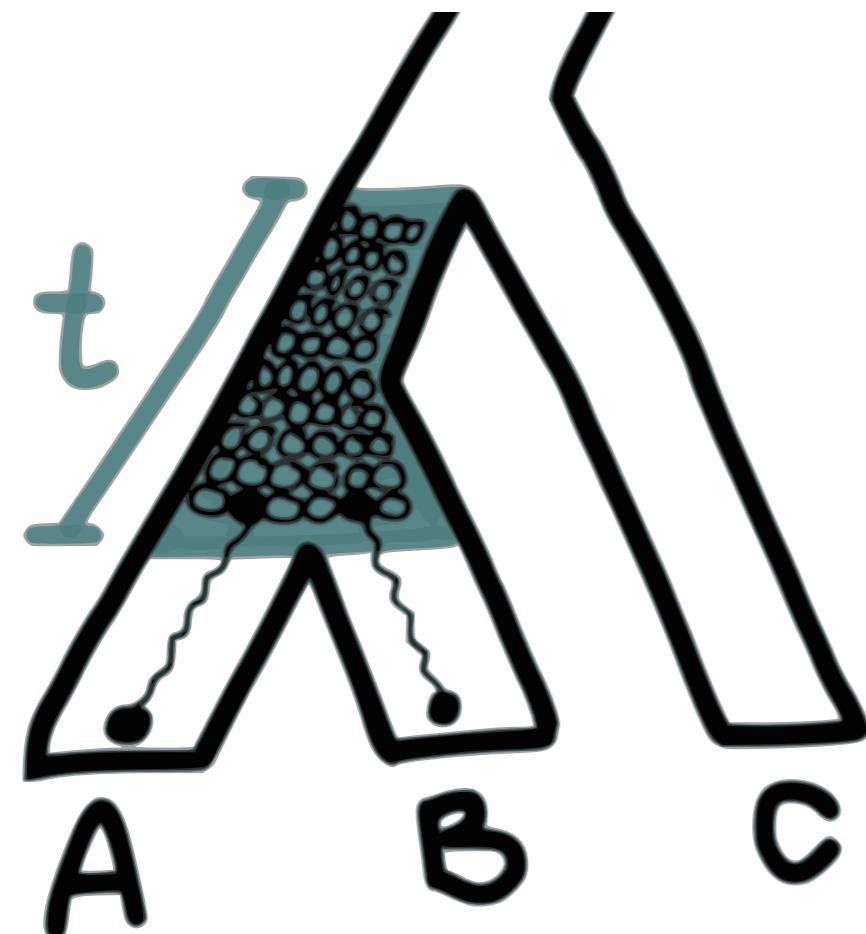
# Multispecies coalescent on a tree



$$P(\wedge_{A B C}) =$$
$$1 - e^{-t}$$
$$+$$

$$P(T > t) = e^{-t}$$

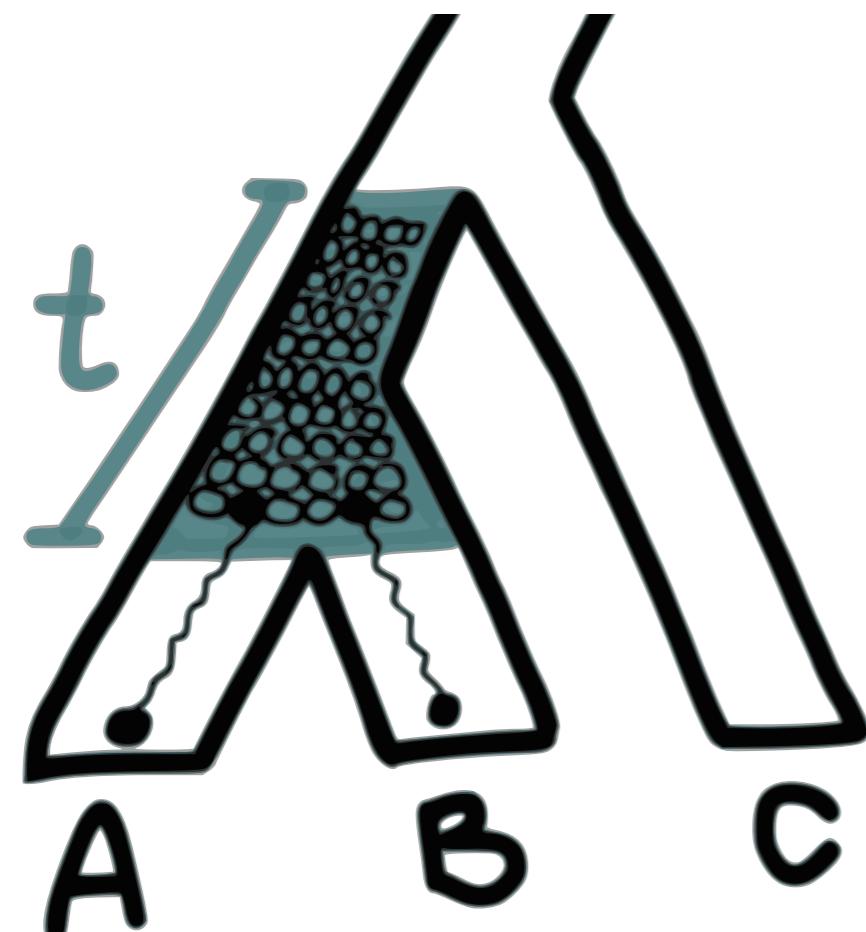
# Multispecies coalescent on a tree



$$P(\wedge_{A B C}) =$$
$$1 - e^{-t}$$
$$+$$
$$e^{-t} \times 1/3$$

$$P(T > t) = e^{-t}$$

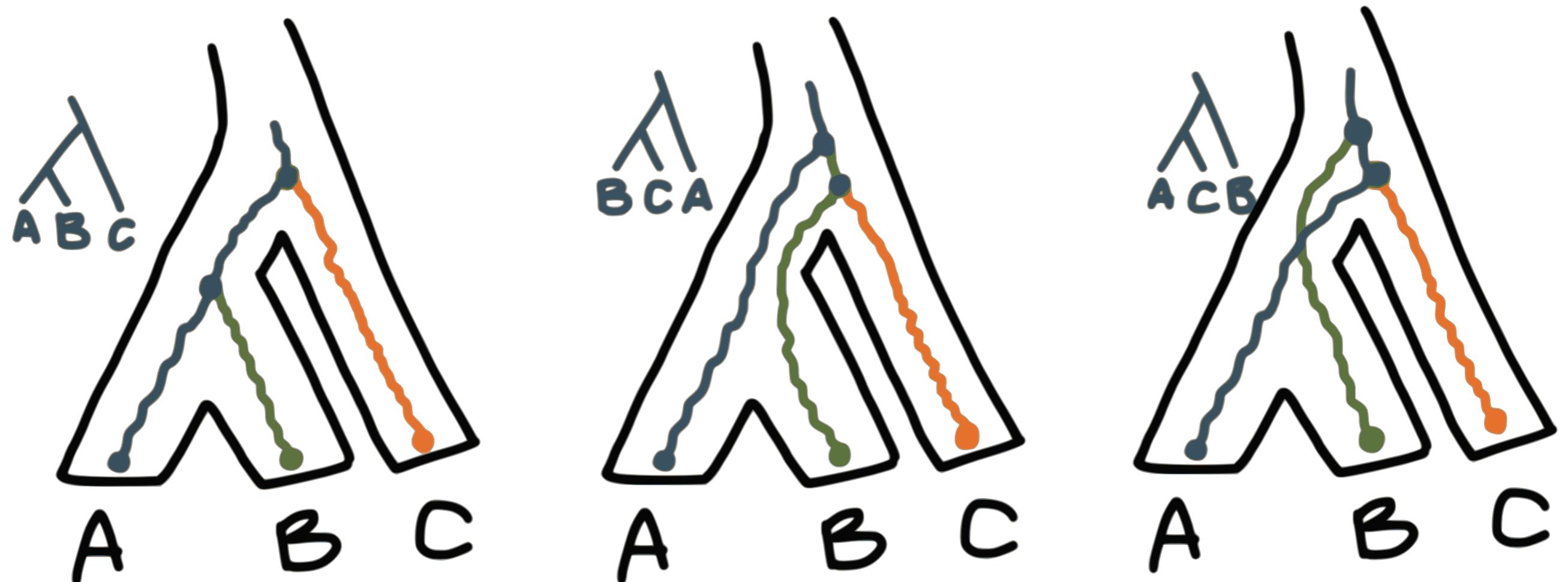
# Multispecies coalescent on a tree



$$P(T > t) = e^{-t}$$

$$\begin{aligned}
& P(\bigwedge_{A \in \mathcal{B}} A) = \\
& 1 - e^{-t} \\
& + \\
& e^{-t} \times 1/3 \\
& = 1 - \frac{2}{3}e^{-t}
\end{aligned}$$

# Multispecies coalescent on a tree



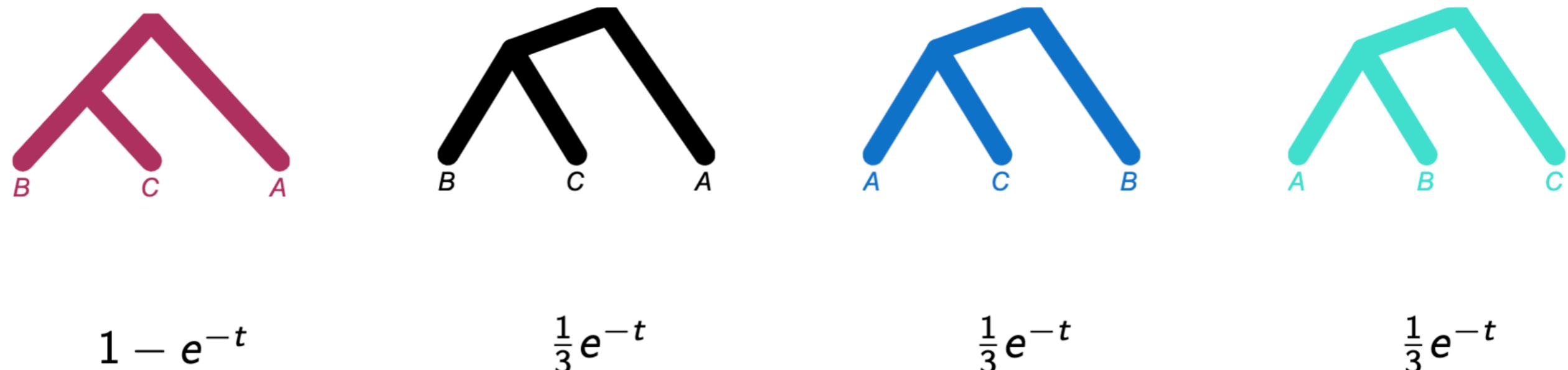
$$1 - \frac{2}{3}e^{-t}$$

$$\frac{1}{3}e^{-t}$$

$$\frac{1}{3}e^{-t}$$

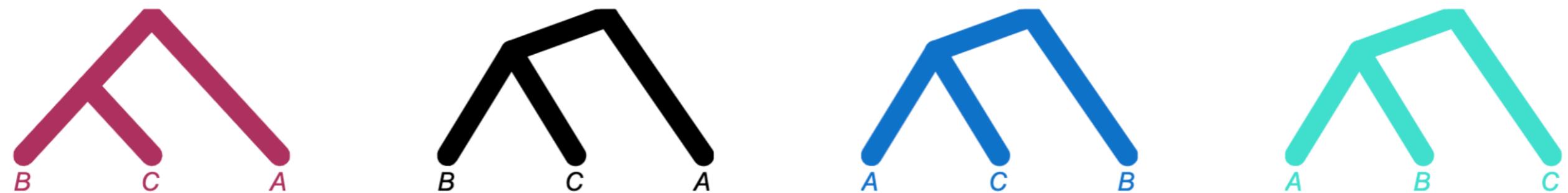
## Phylogenetic coalescent model

Probabilities of each gene tree history are shown below them  
 $t$  = length of interval between speciation events



## Phylogenetic coalescent model

$t = \text{length of interval between coalescent events} = 1.0 = 0.5 = 2.0$



$$1 - e^{-t}$$

0.63

0.40

0.85

$$\frac{1}{3}e^{-t}$$

0.12

0.20

0.05

$$\frac{1}{3}e^{-t}$$

0.12

0.20

0.05

$$\frac{1}{3}e^{-t}$$

0.12

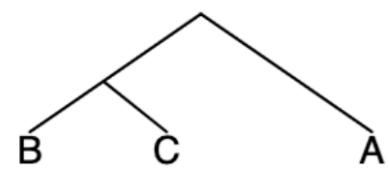
0.20

0.05

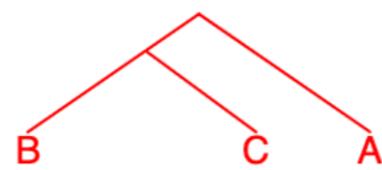
## Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

- What are these probabilities like as a function of  $t$ , the length of time between speciation events?

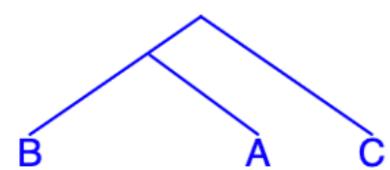
(b)



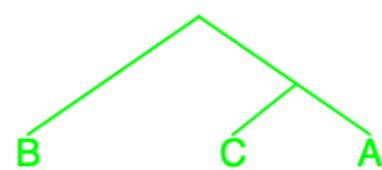
$$\text{prob} = 1 - \exp(-t)$$



$$\text{prob} = (1/3)\exp(-t)$$

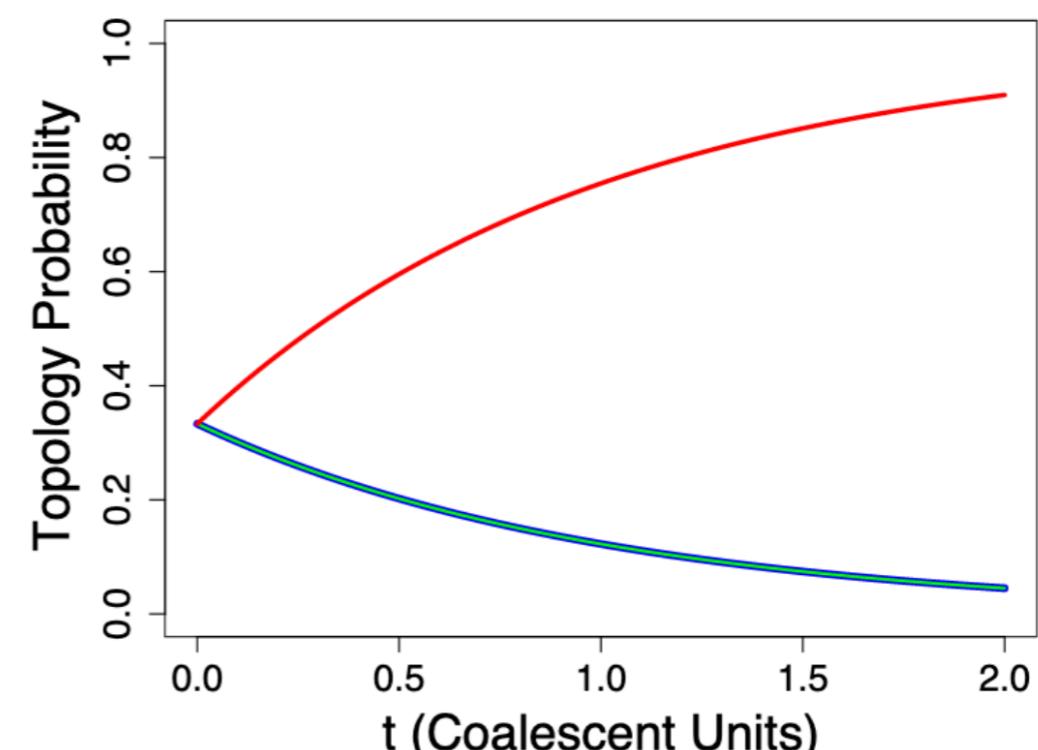


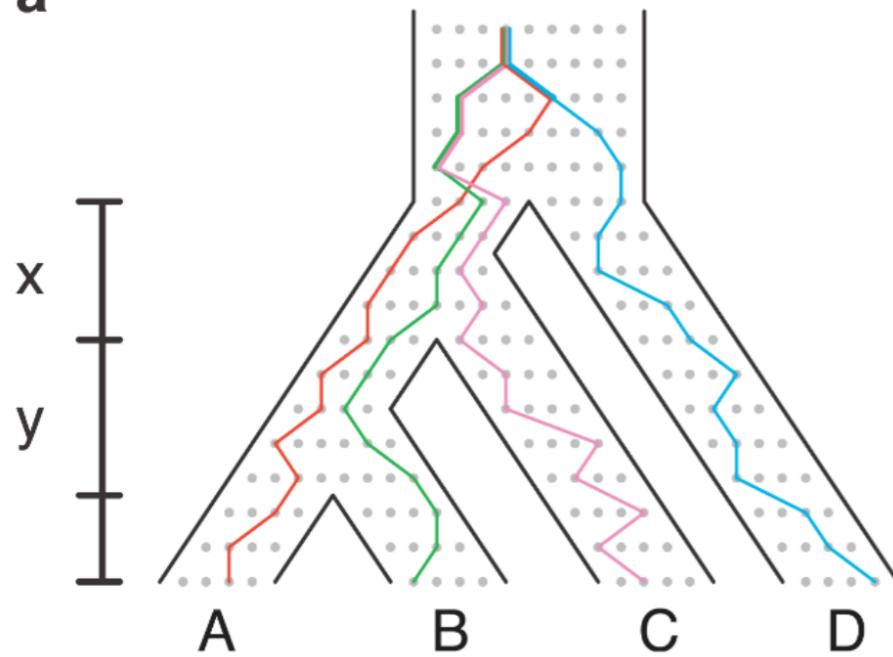
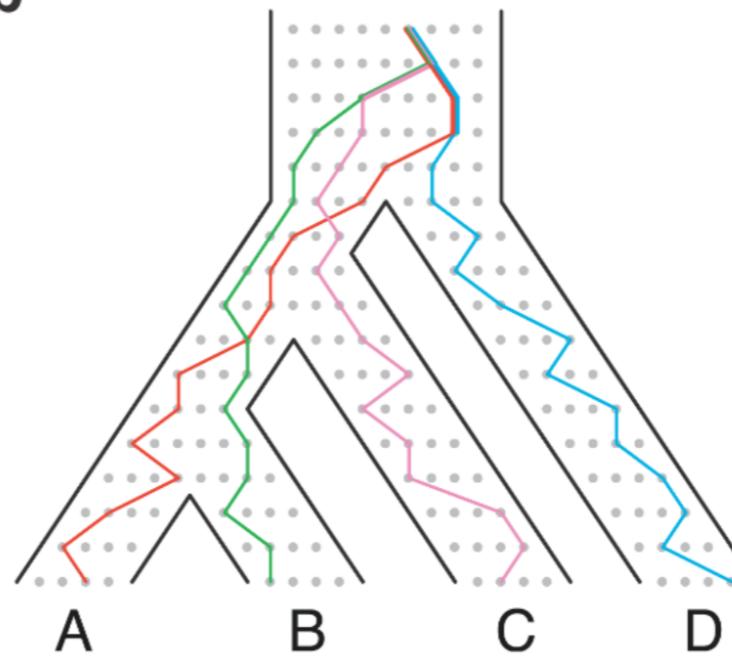
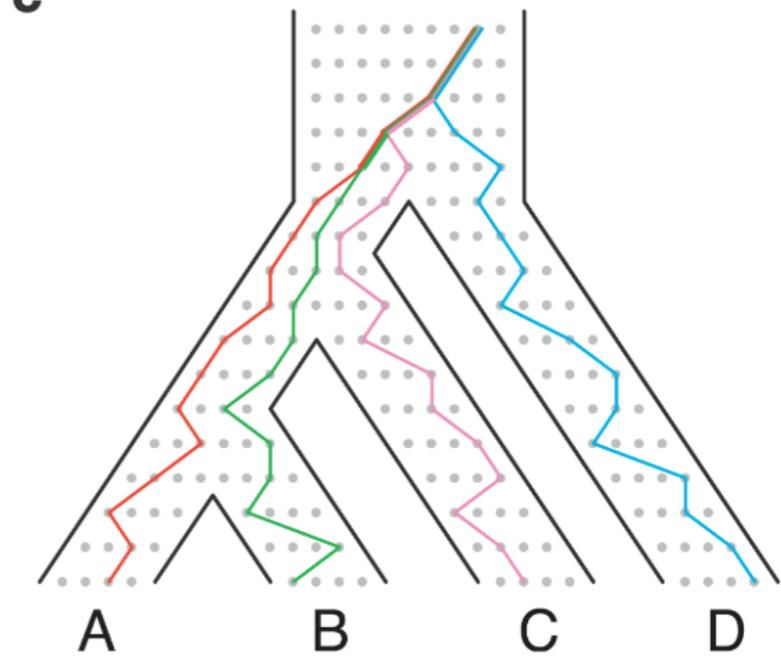
$$\text{prob} = (1/3)\exp(-t)$$



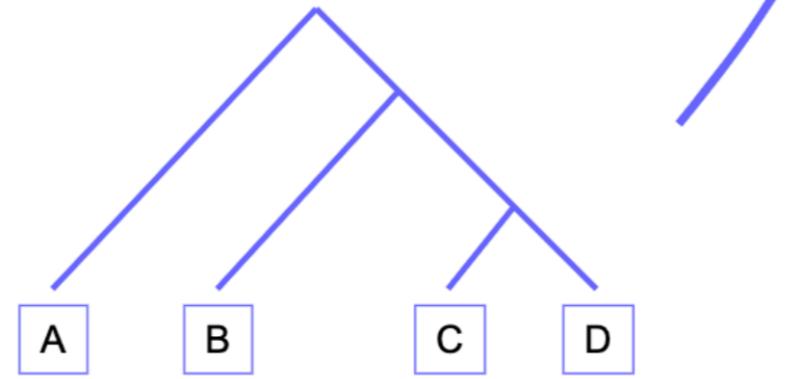
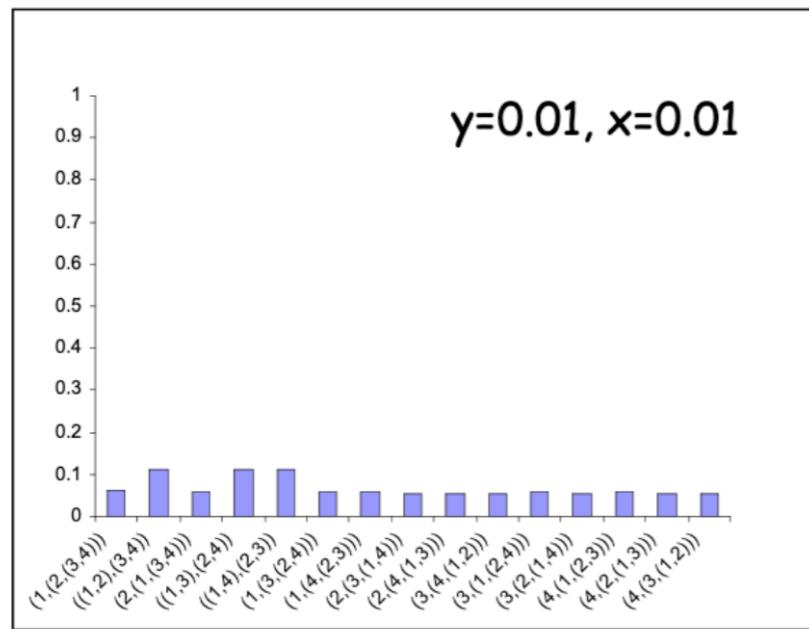
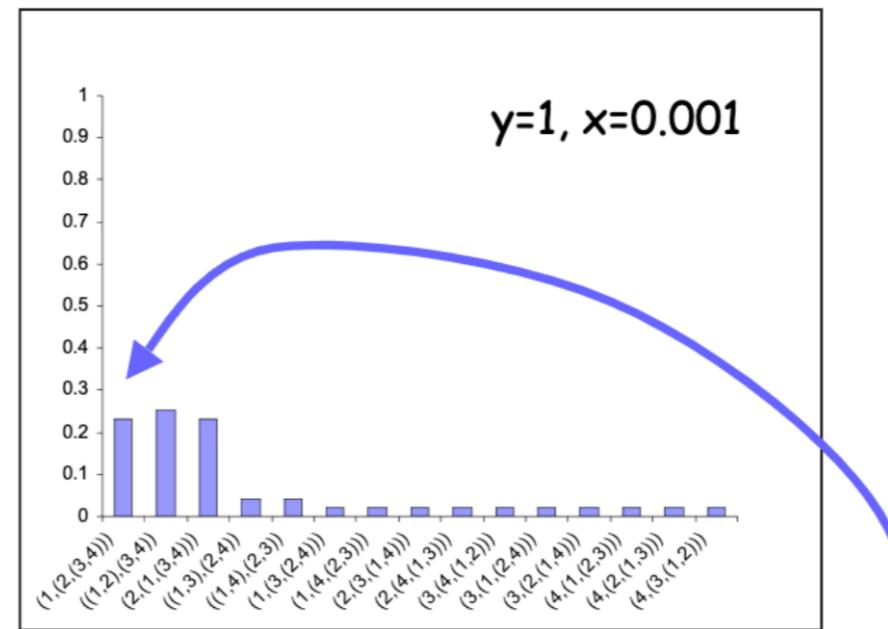
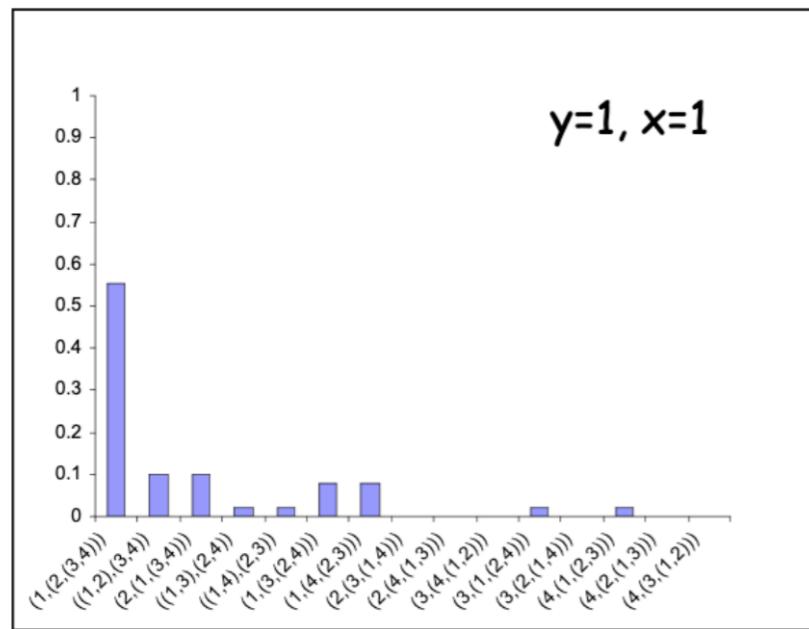
$$\text{prob} = (1/3)\exp(-t)$$

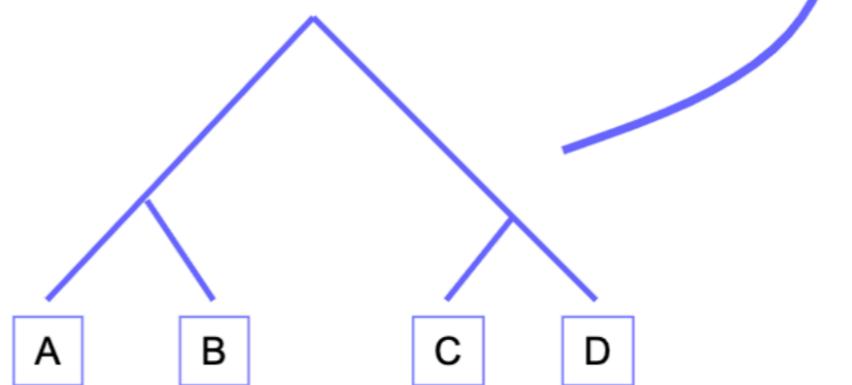
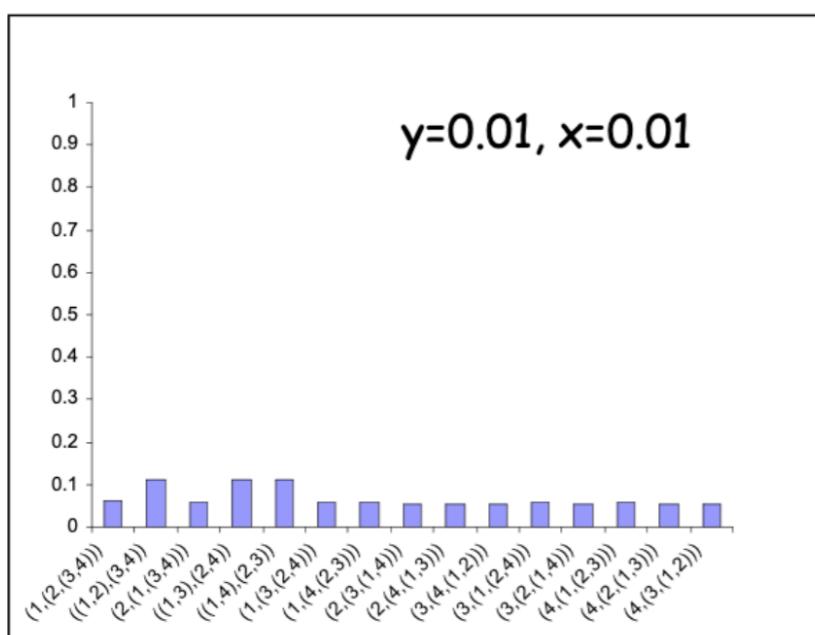
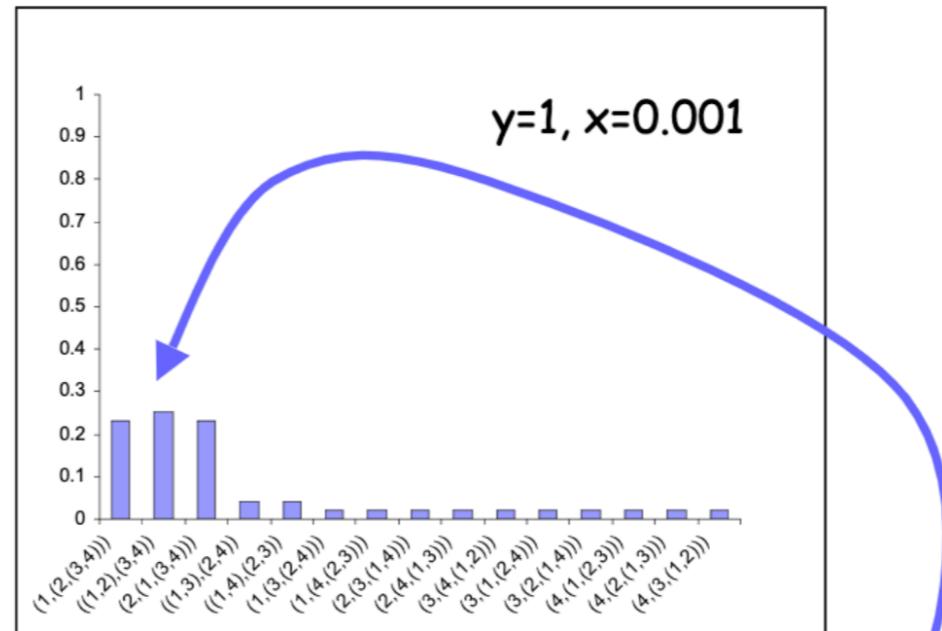
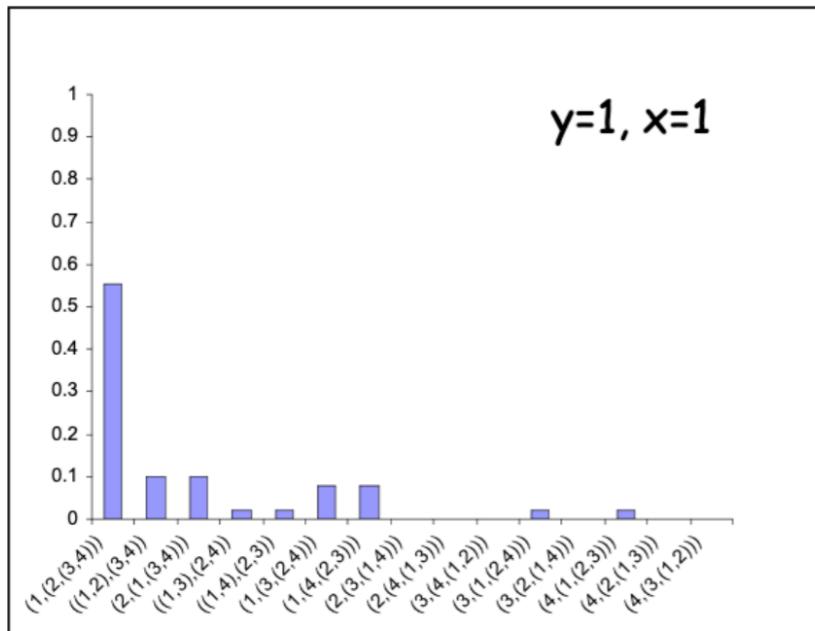
(c)

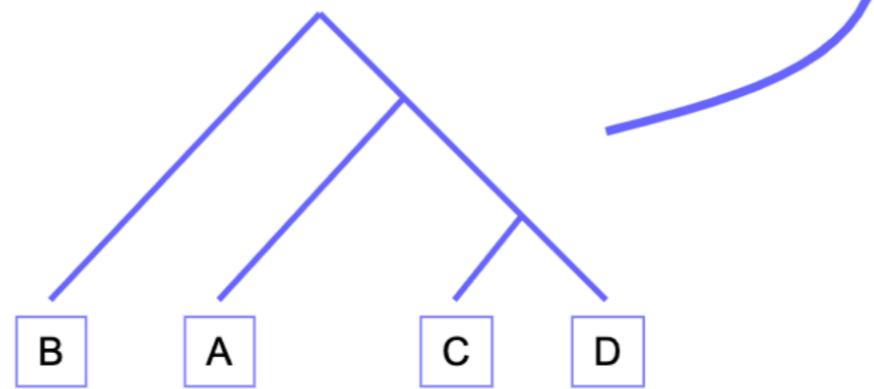
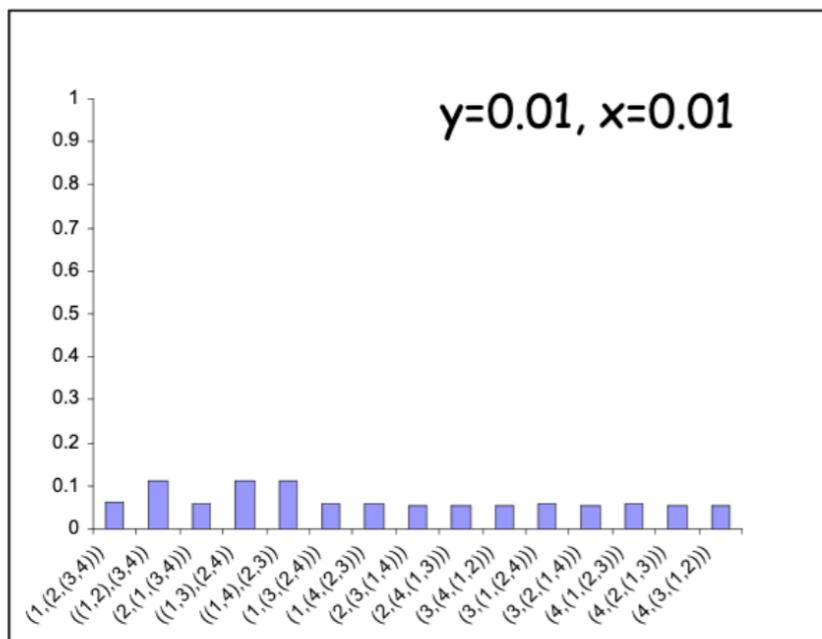
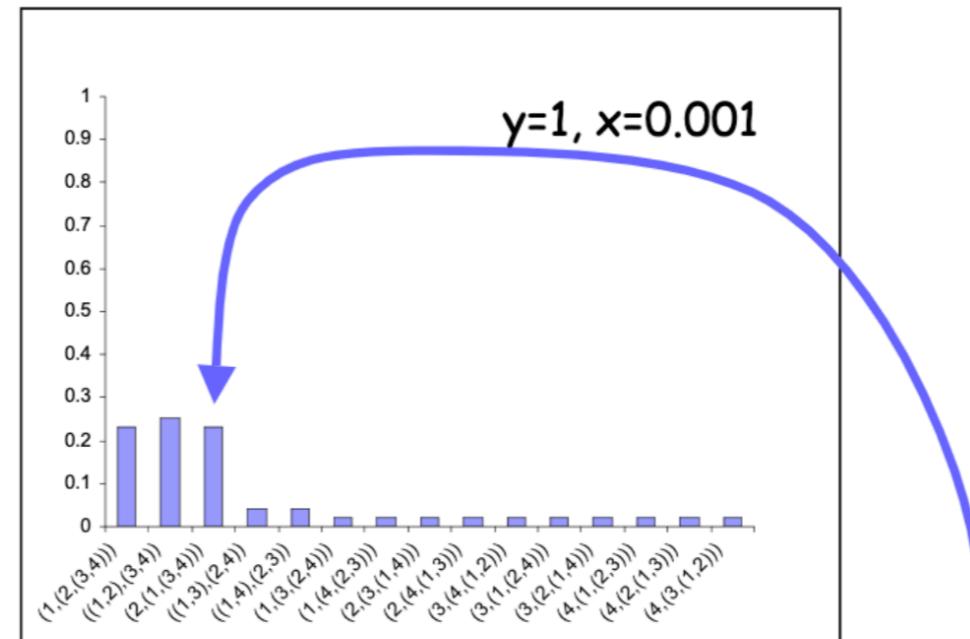
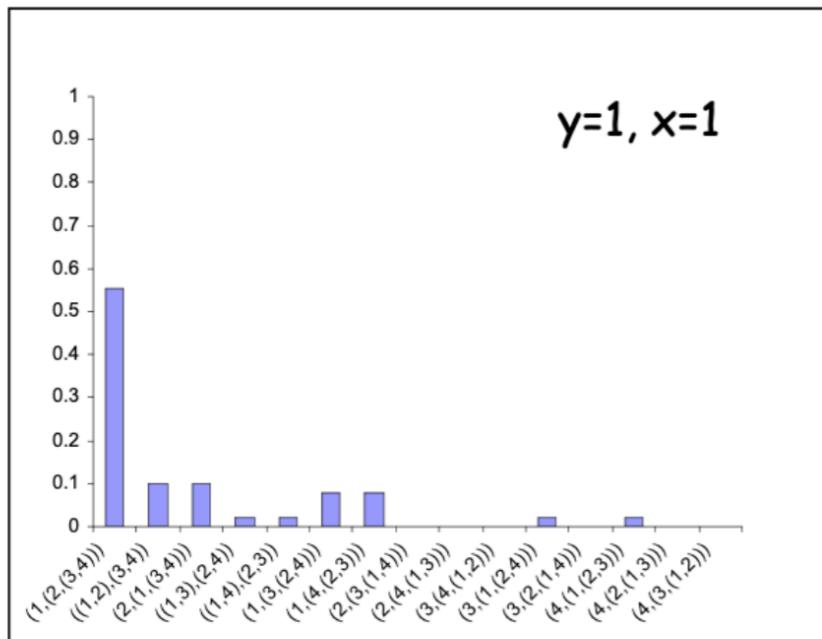


**a****b****c**

If the internal branches of the species tree—*x* and *y*—are short so that coalescences occur deep in the tree, the two sequences of coalescences that produce a given symmetric gene tree topology together have higher probability than the single sequence that produces the topology that matches the species tree. (a) and (b) Two coalescence sequences leading to gene tree topology ((AD)(BC)). In (a), the lineages from B and C coalesce more recently than those from A and D, and in (b), the reverse is true. (c) The single sequence of coalescences leading to gene tree topology (((AB)C)D). (Degnan and Rosenberg, 2006)







# Anomaly zone

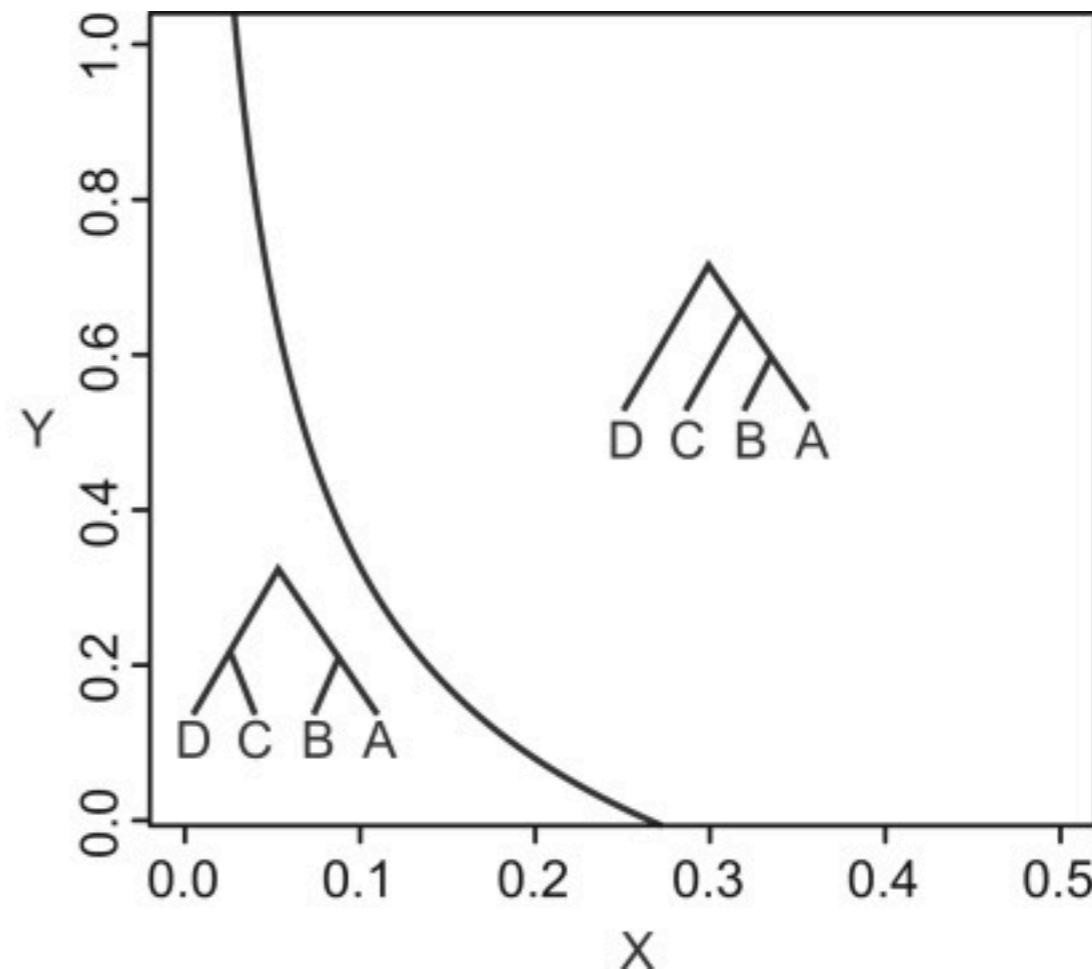
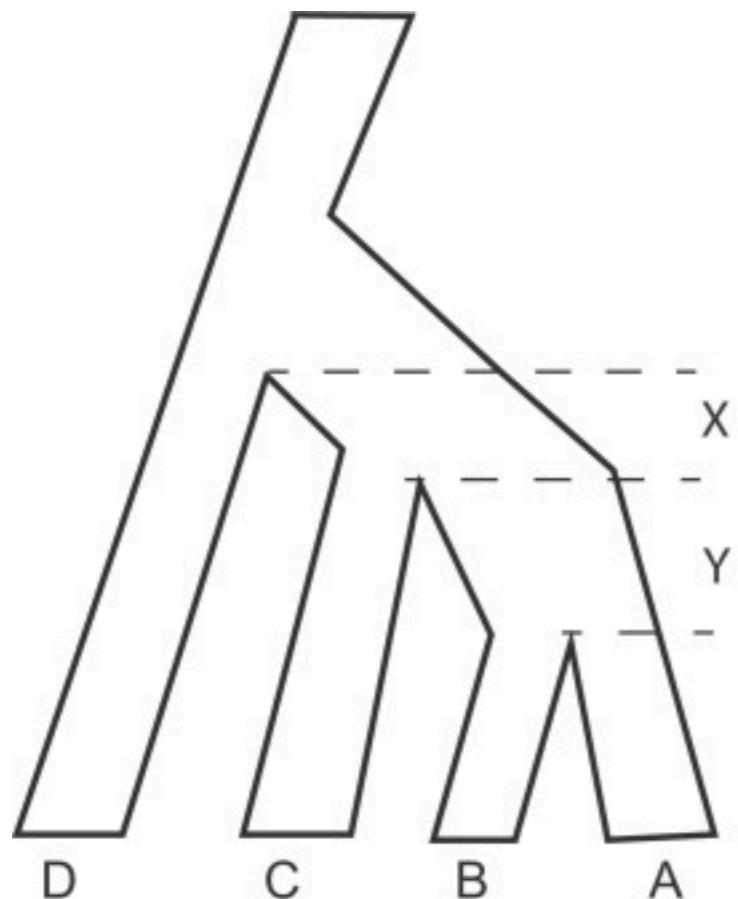


Figure 1 in [Linkem et al \(2016\)](#)



Distances  
Parsimony  
Likelihood  
(Bayesian)

$$P(T, \theta | G) \propto \pi(T)\pi(\theta) \prod_{i=1}^L P(G_i | T, \theta)$$

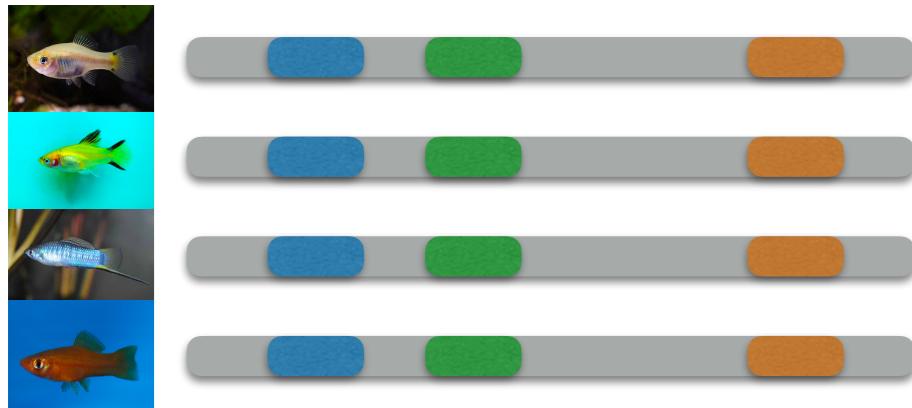
**Prior Tree**      **Multispecies Coalescent**      **Bayesian**

$$L(T, \theta) = \prod_{i=1}^L P(G_i | T, \theta)$$

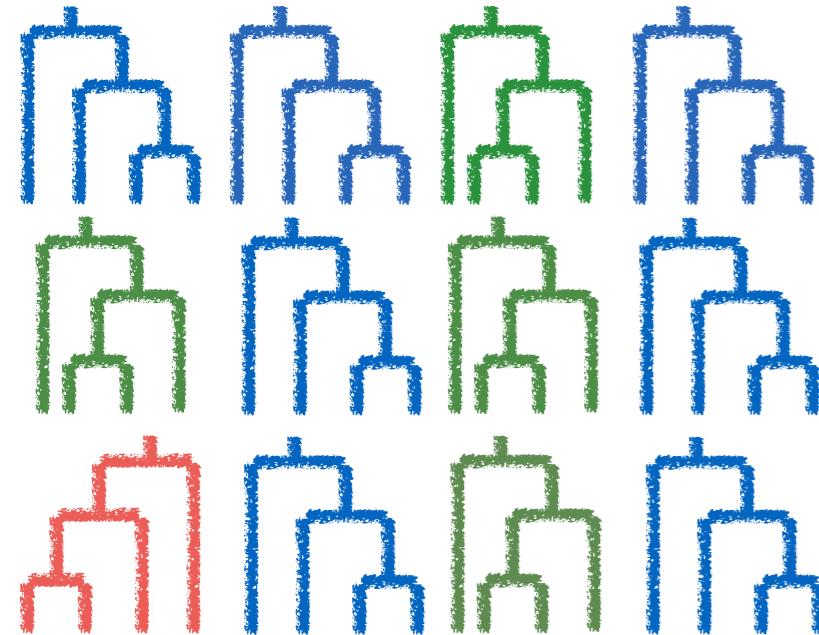
**Max. Lik.**

Summary methods: ASTRAL, BUCKy, MP-EST

(Zhang et al, 2018) (Larget et al, 2010) (Liu et al, 2010)



Distances  
Parsimony  
Likelihood  
(Bayesian)



$$P(T, \theta | G) \propto \pi(T) \pi(\theta) \prod_{i=1}^L P(G_i | T, \theta)$$

Prior Tree      Multispecies Coalescent



$$L(T, \theta) = \prod_{i=1}^L P(G_i | T, \theta)$$

Bayesian

Max. Lik.



Summary methods: ASTRAL, BUCKy, MP-EST

(Zhang et al, 2018) (Larget et al, 2010) (Liu et al, 2010)



Distances  
Parsimony  
Likelihood  
(Bayesian)

$$P(T, \theta | G) \propto \pi(T)\pi(\theta) \prod_{i=1}^L P(G_i | T, \theta)$$

**Prior Tree**      **Multispecies Coalescent**      **Bayesian**

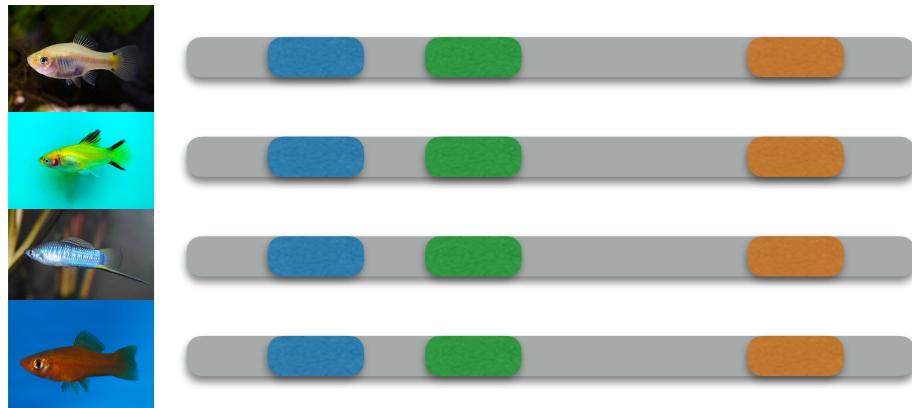
$$L(T, \theta) = \prod_{i=1}^L P(G_i | T, \theta)$$

**Max. Lik.**

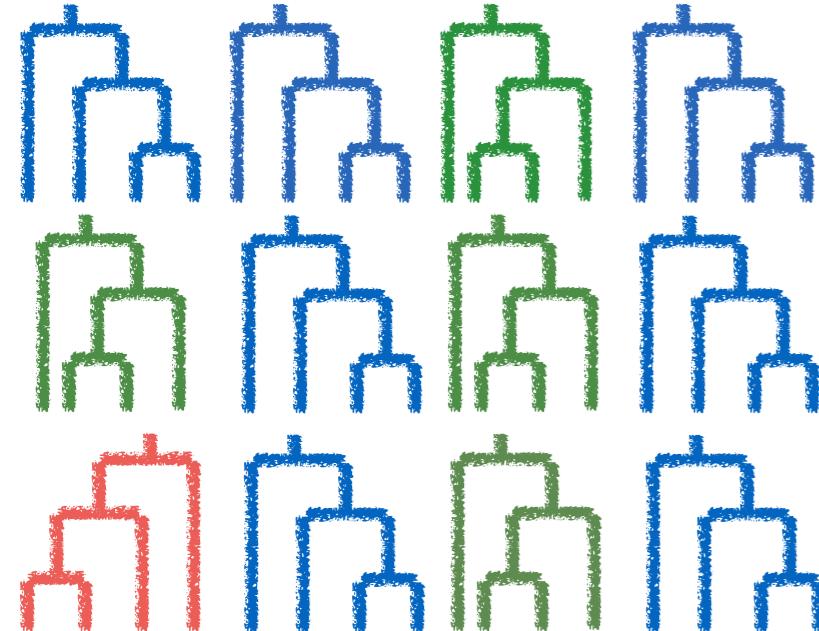
Do not search tree space

Summary methods: ASTRAL, BUCKy, MP-EST

(Zhang et al, 2018) (Larget et al, 2010) (Liu et al, 2010)



Distances  
Parsimony  
Likelihood  
(Bayesian)



$$P(T, \theta | G) \propto \pi(T) \pi(\theta) \prod_{i=1}^L P(G_i | T, \theta)$$

Prior Tree      Multispecies Coalescent



$$L(T, \theta) = \prod_{i=1}^L P(G_i | T, \theta)$$

Bayesian

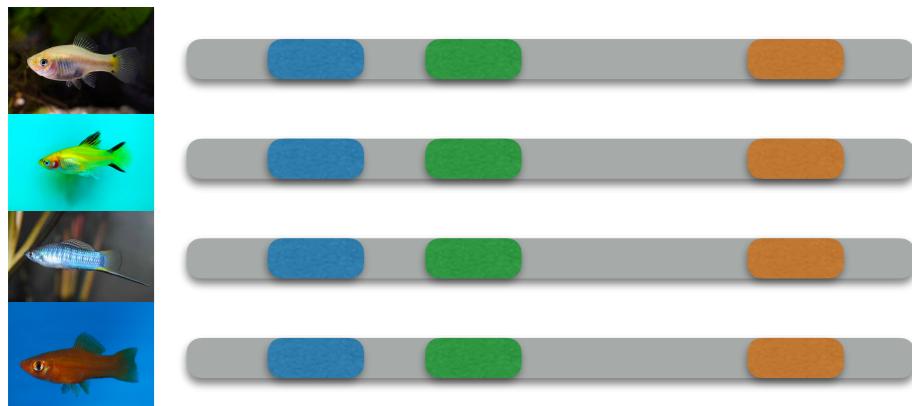
Max. Lik.

Do not search tree space

Summary methods: ASTRAL, BUCKy, MP-EST

(Zhang et al, 2018) (Larget et al, 2010) (Liu et al, 2010)

## Co-estimation (lecture 15)



$$P(T, G, \theta | D) \propto \pi(T)\pi(\theta) \prod_{i=1}^L P(D_i | G_i)P(G_i | T)$$

Tree  
 ↓  
 $P(T, G, \theta | D)$   
 ↑ gene trees    ↑ L loci    ↑ Prior Tree  
 Substitution Model    Multispecies Coalescent

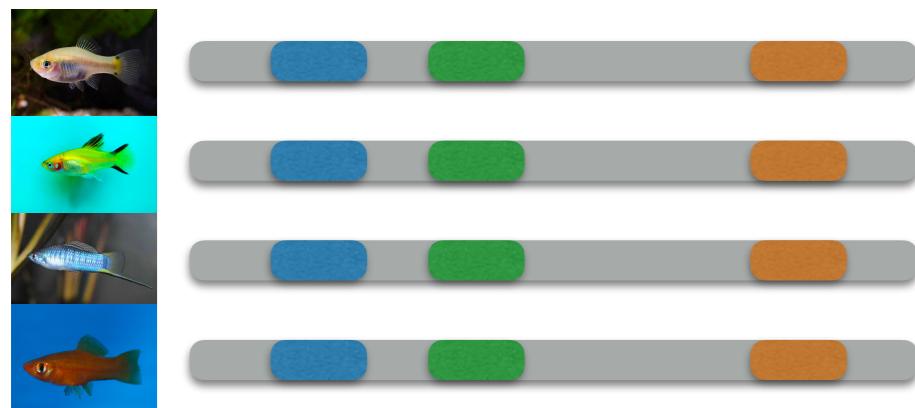
$$P(T, \theta | G) \propto \pi(T)\pi(\theta) \prod_{i=1}^L P(G_i | T, \theta) \quad \text{Bayesian}$$

$$L(T, \theta) = \prod_{i=1}^L P(G_i | T, \theta) \quad \text{Max. Lik.}$$

Summary methods: ASTRAL, BUCKy, MP-EST

(Zhang et al, 2018) (Larget et al, 2010) (Liu et al, 2010)

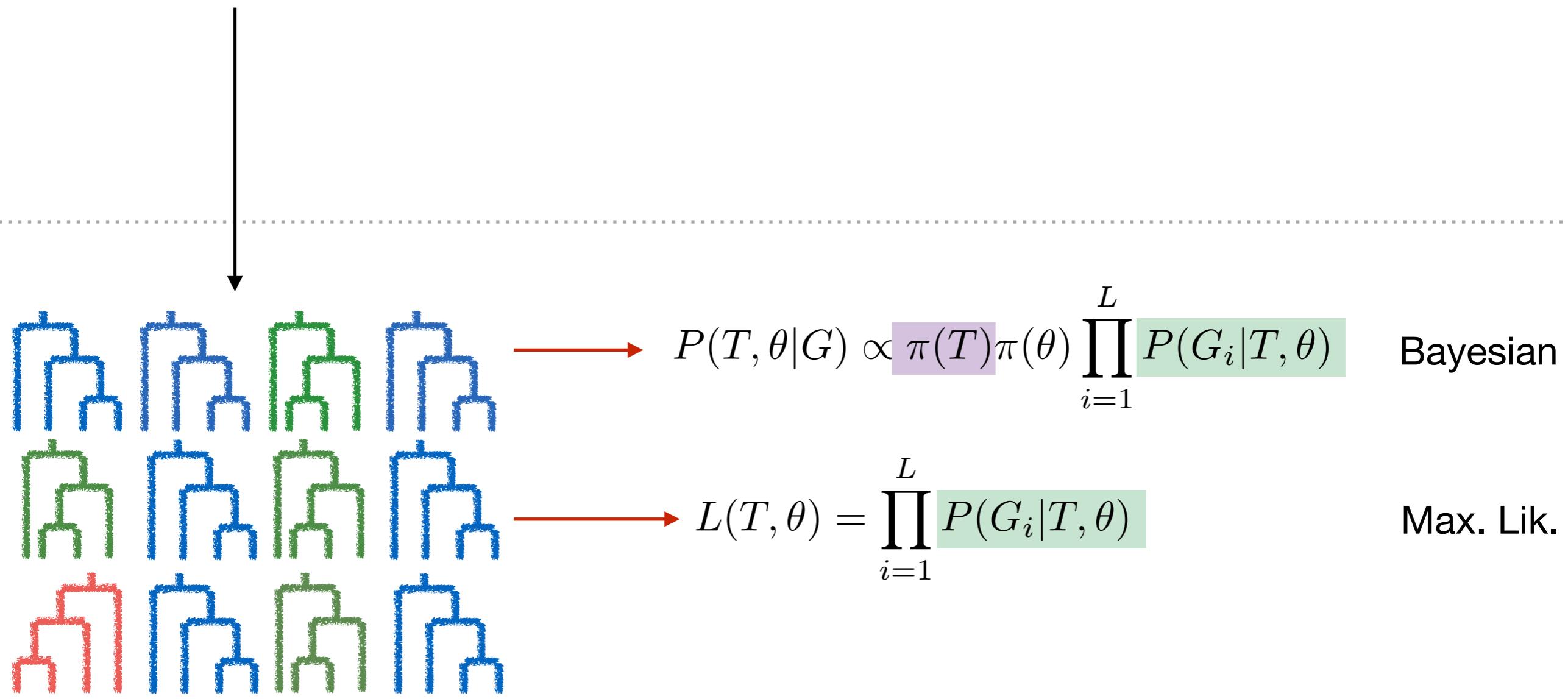
## Co-estimation (lecture 15)



$$P(T, G, \theta | D) \propto \pi(T)\pi(\theta) \prod_{i=1}^L P(D_i | G_i)P(G_i | T)$$

Tree  
gene trees      L loci  
↑                ↑  
**Prior Tree**

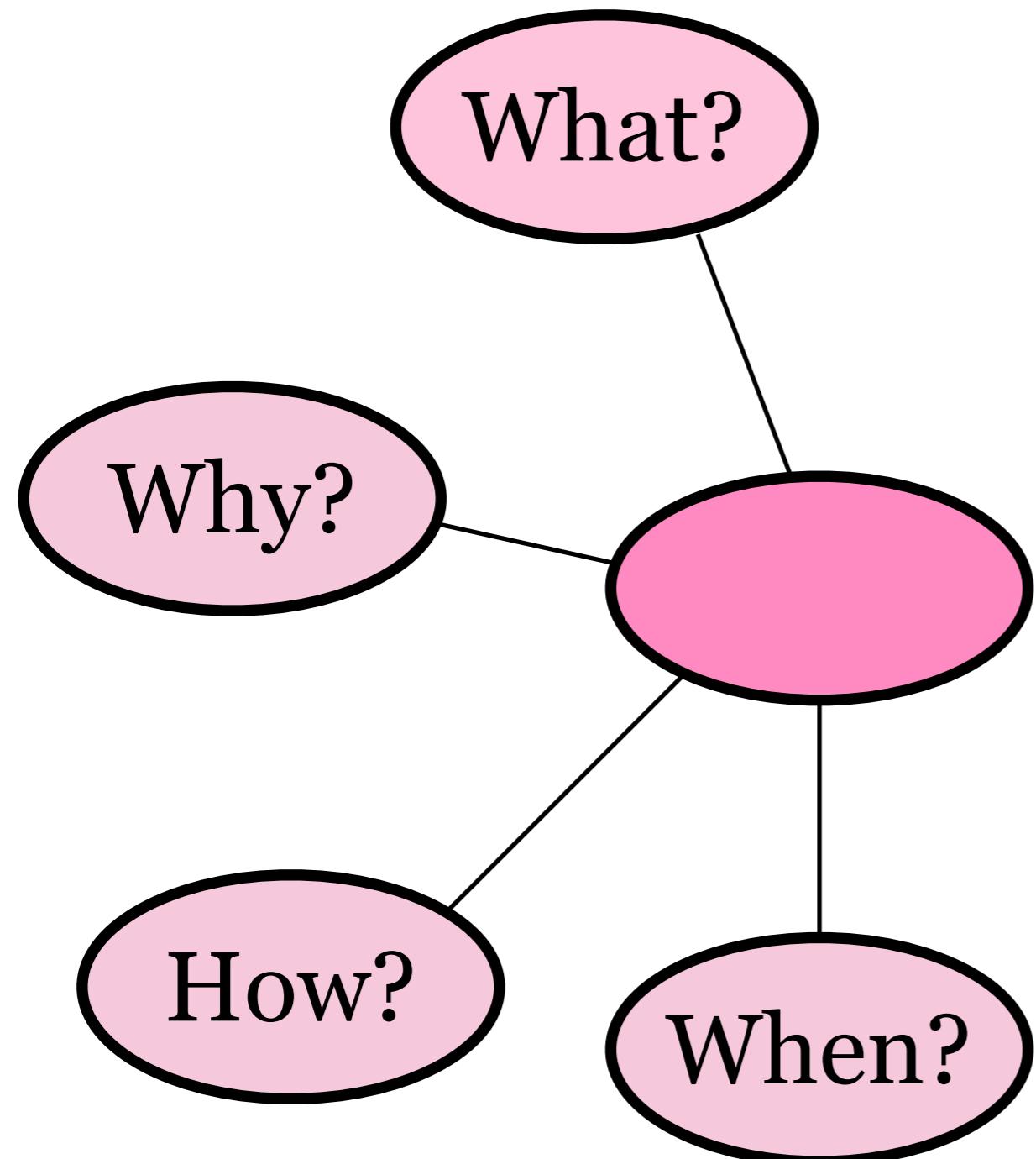
**Substitution Model**    **Multispecies Coalescent**



Summary methods: ASTRAL, BUCKy, MP-EST

(Zhang et al, 2018) (Larget et al, 2010) (Liu et al, 2010)

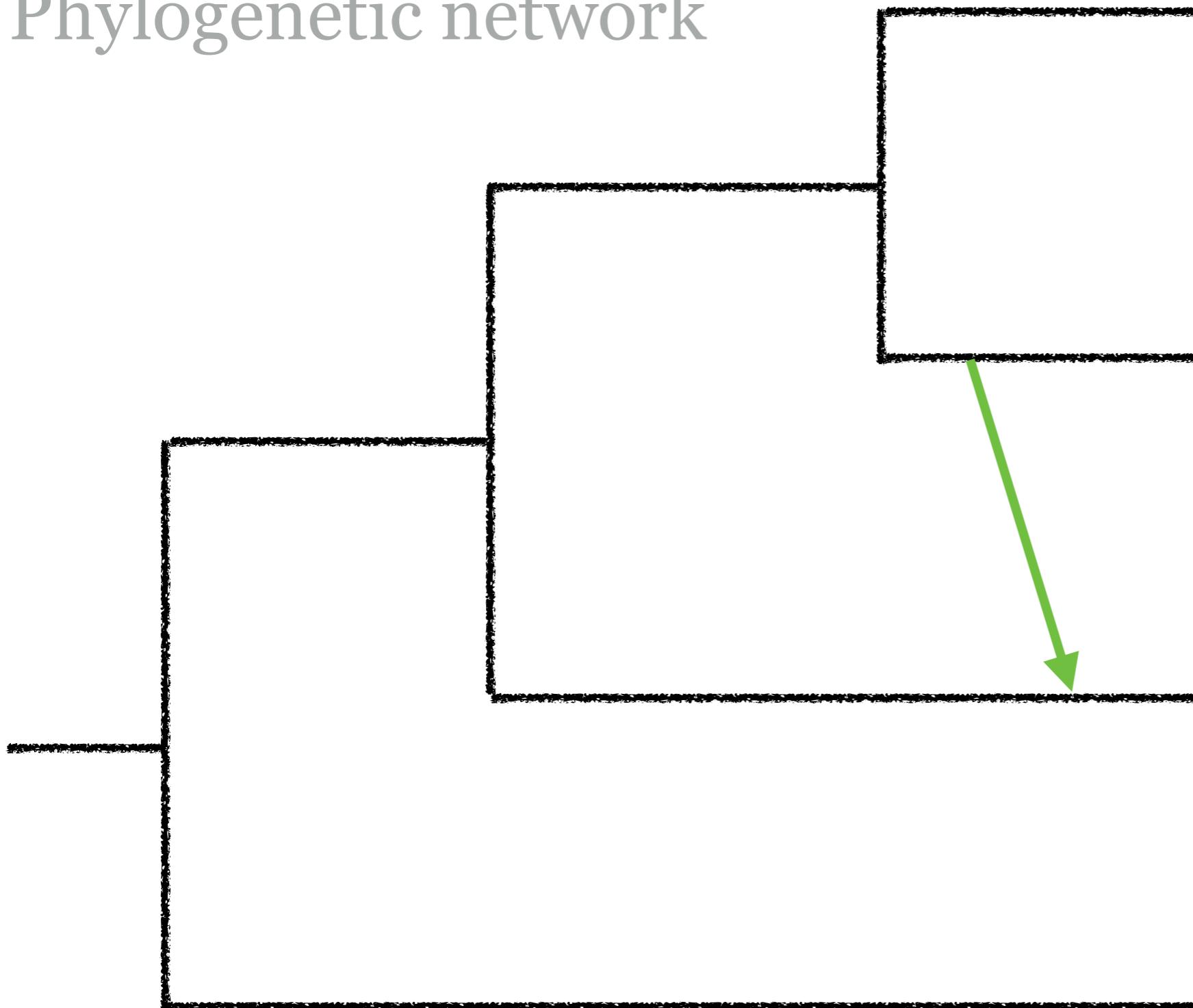
- Concatenation: assumes all genes follow the same tree-like history
- Coalescent-based tree methods: accounts for ILS (and sometimes gene tree estimation error), but does not allow for other sources of gene tree discordance (like GDL or gene flow)
- Extensions to coalescent-based tree methods:
  - Coalescent-based network methods: accounts for ILS, gene tree estimation error and gene flow
  - Coalescent+GDL methods (Li et al, 2020)



# Phylogenetic Networks

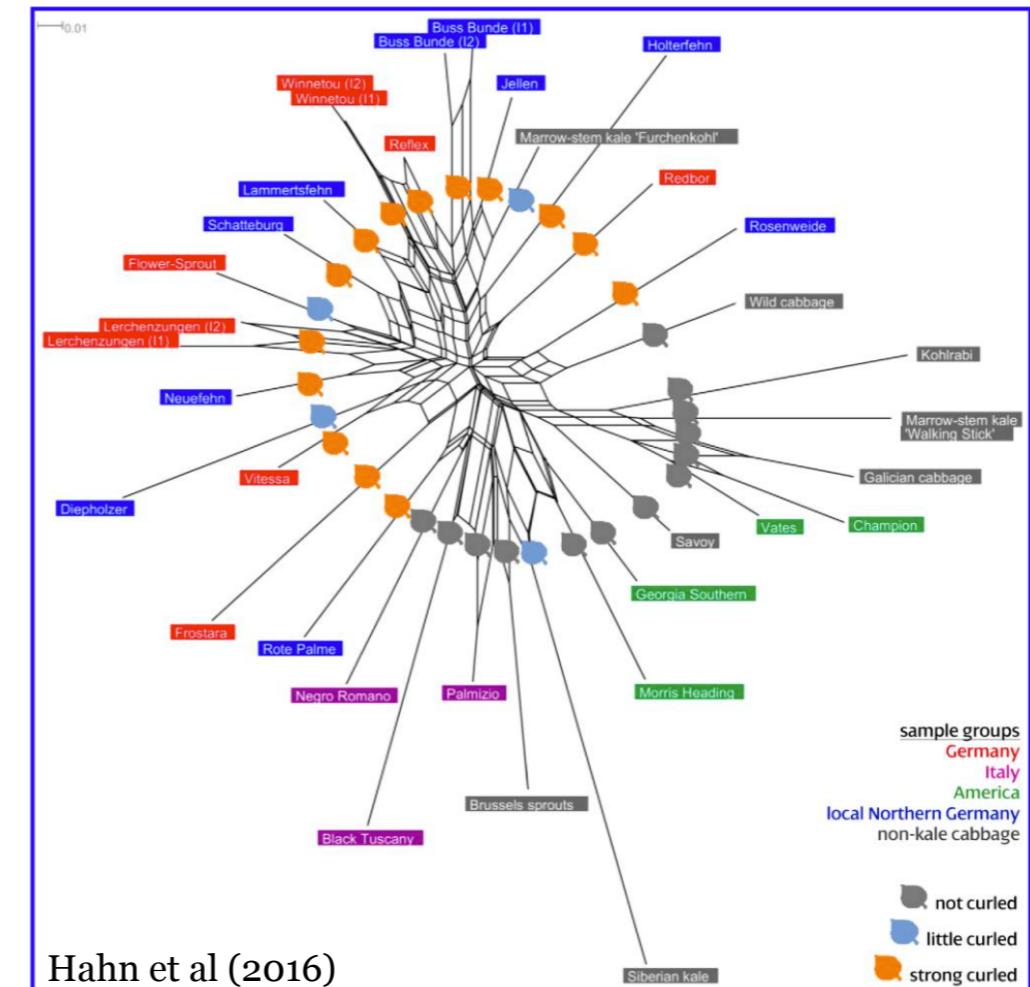
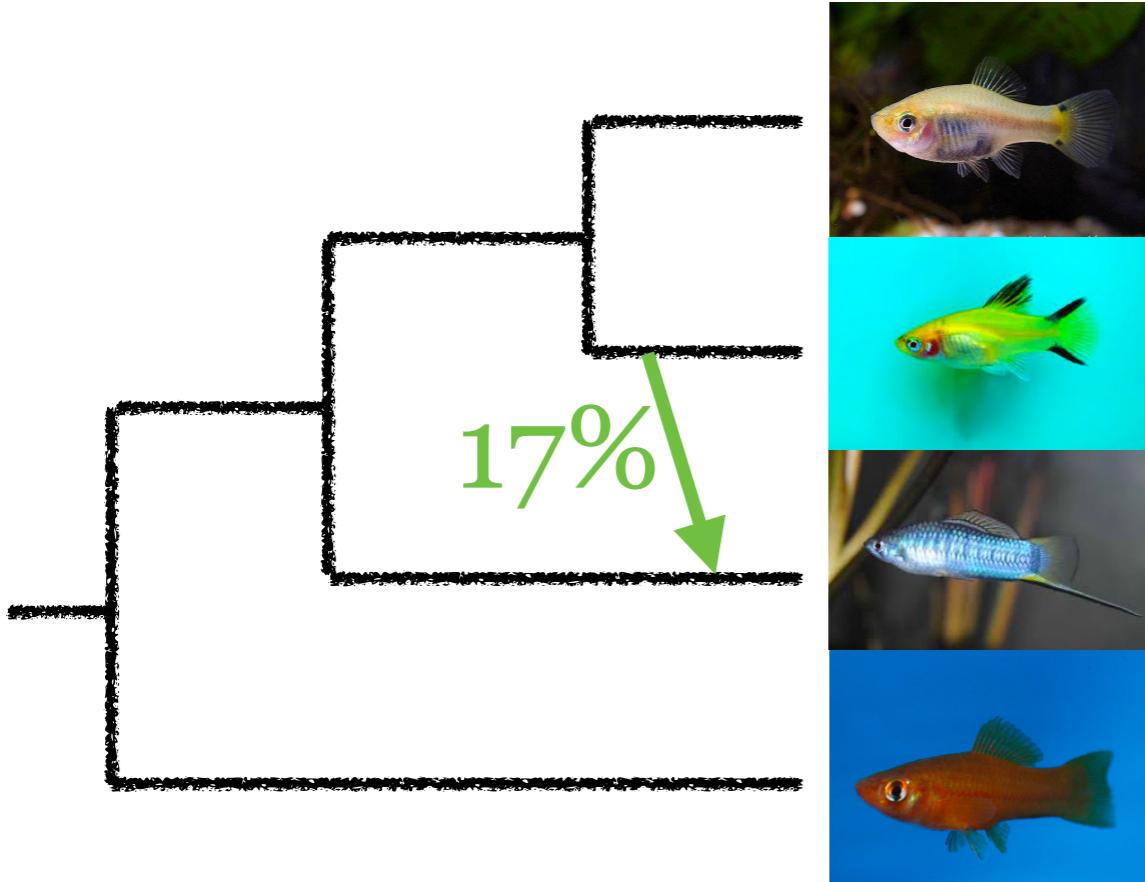
# What?

Phylogenetic network



# What?

## Phylogenetic network

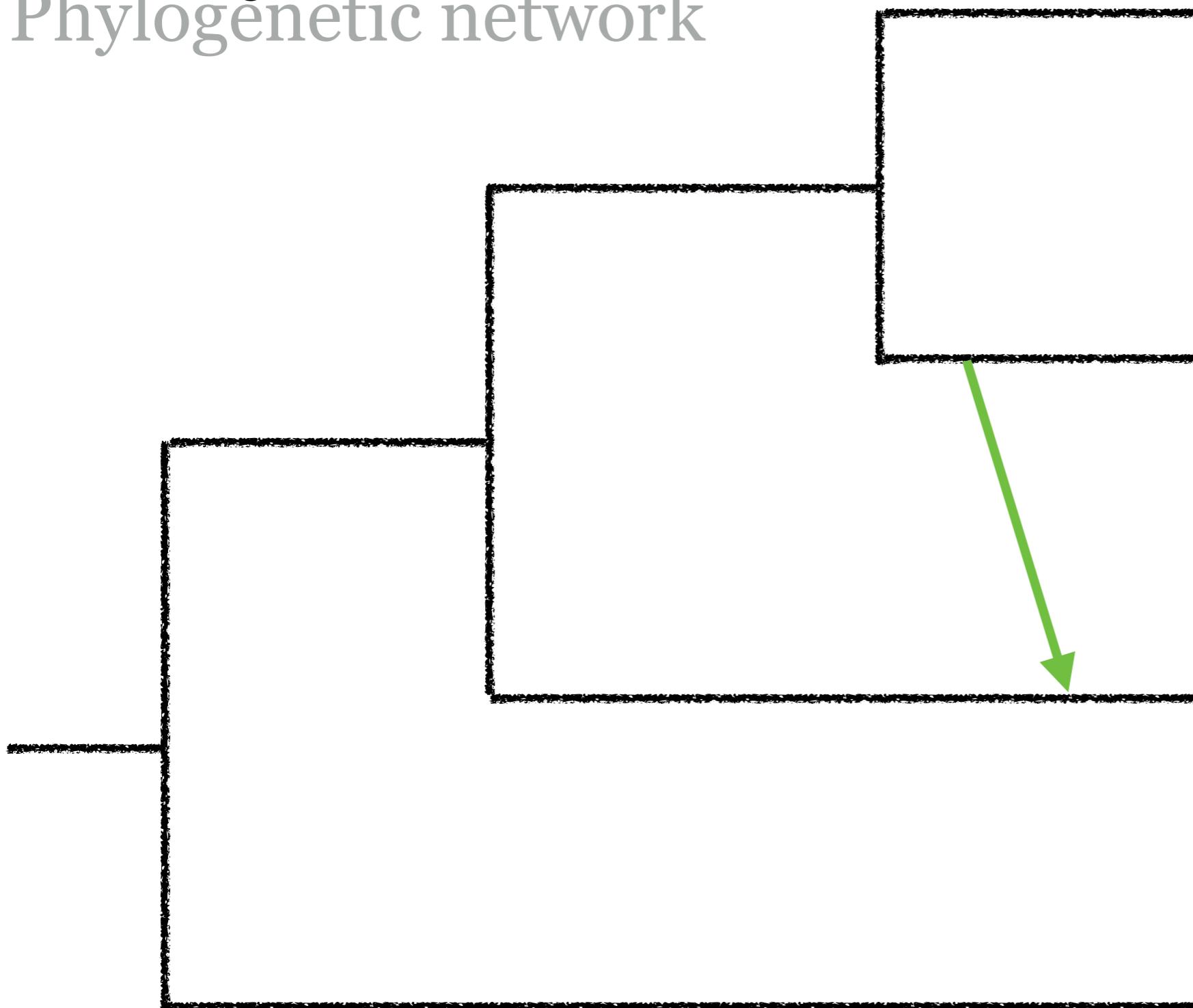


Explicit

Implicit

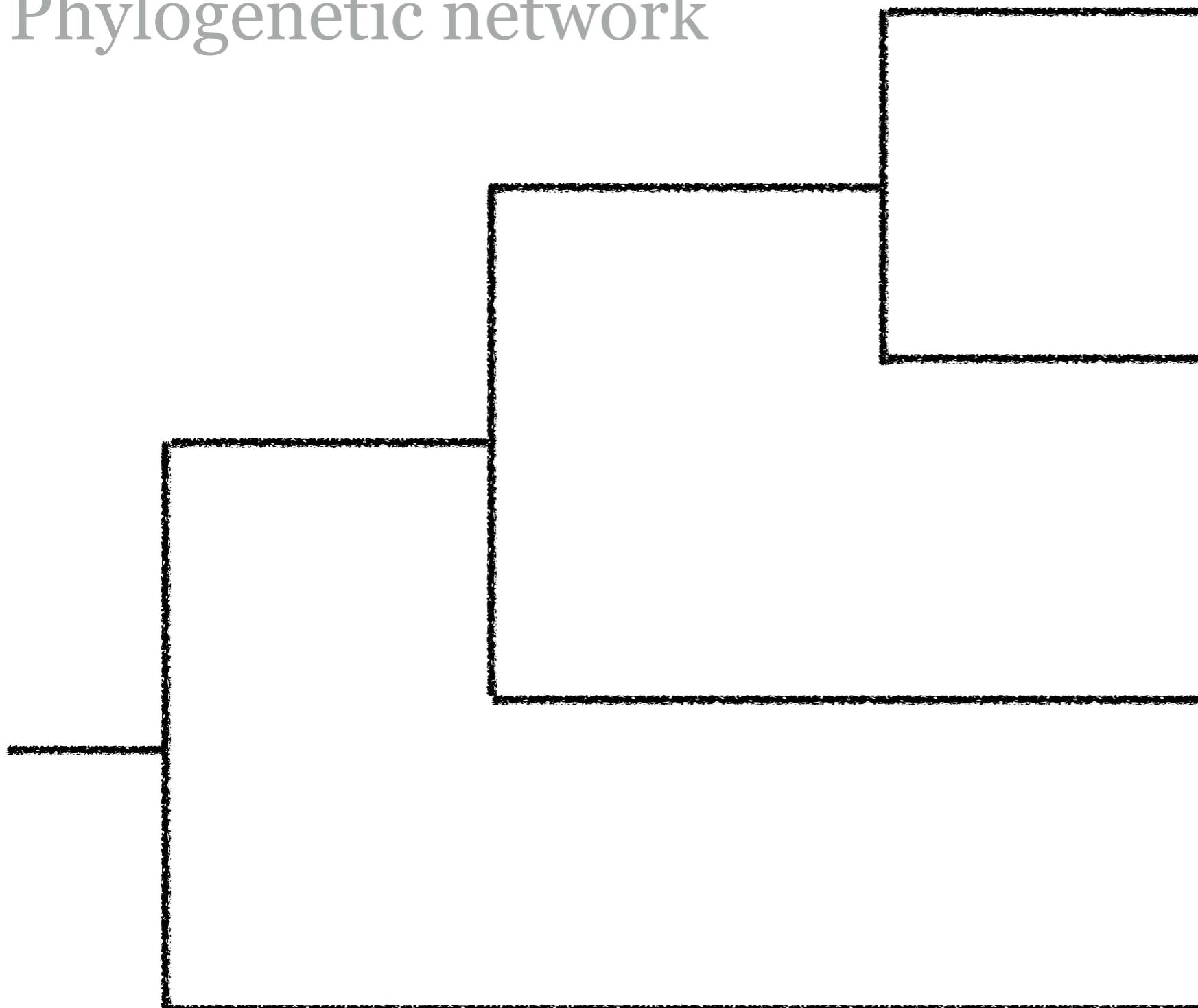
# Why?

# Phylogenetic network



# Why?

Phylogenetic network



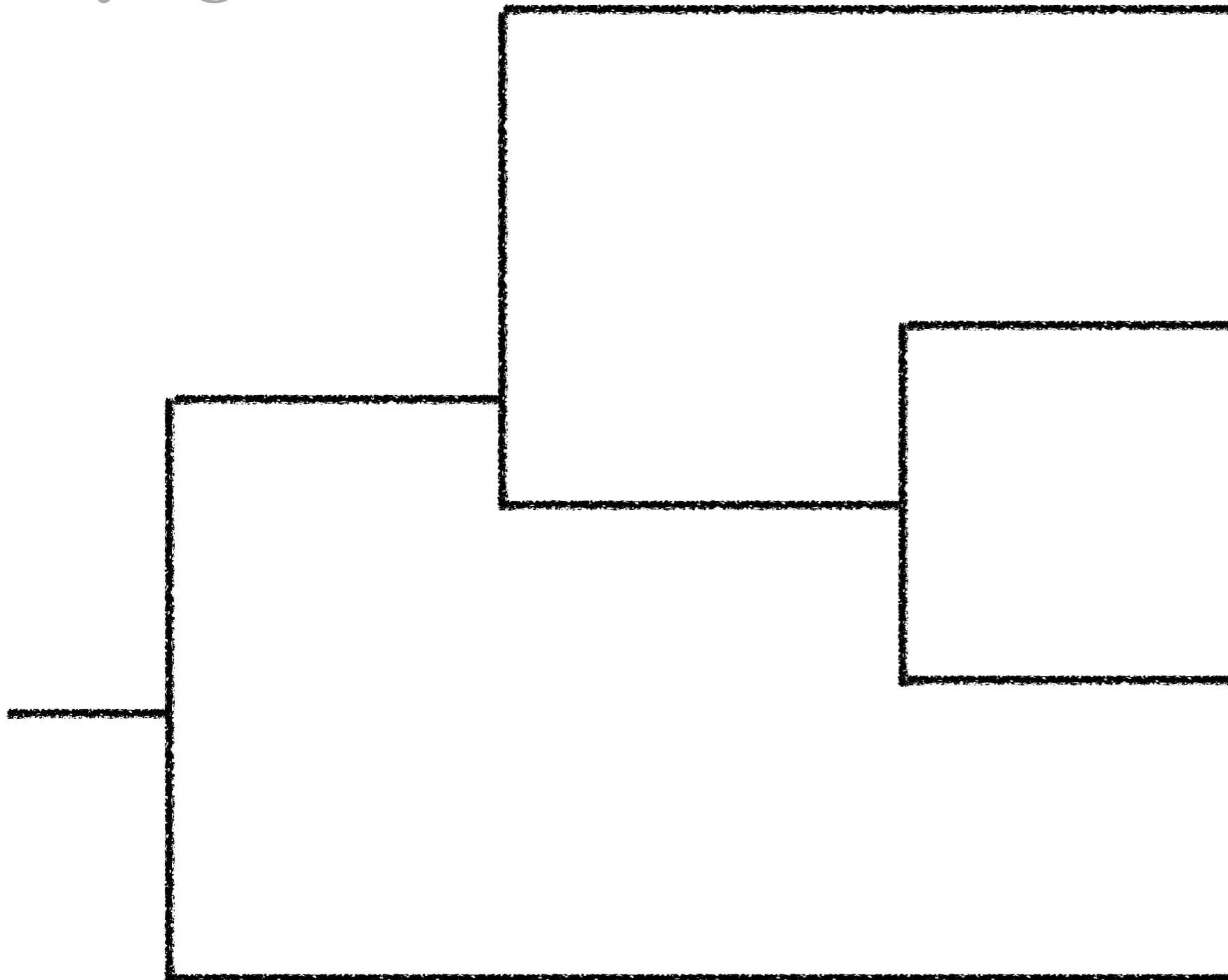
Main tree



# Why?

Phylogenetic network

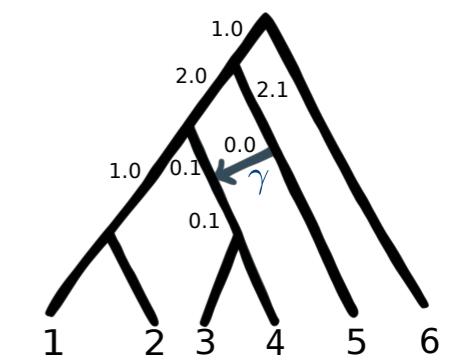
Ignore gene flow  
=>Wrong tree!



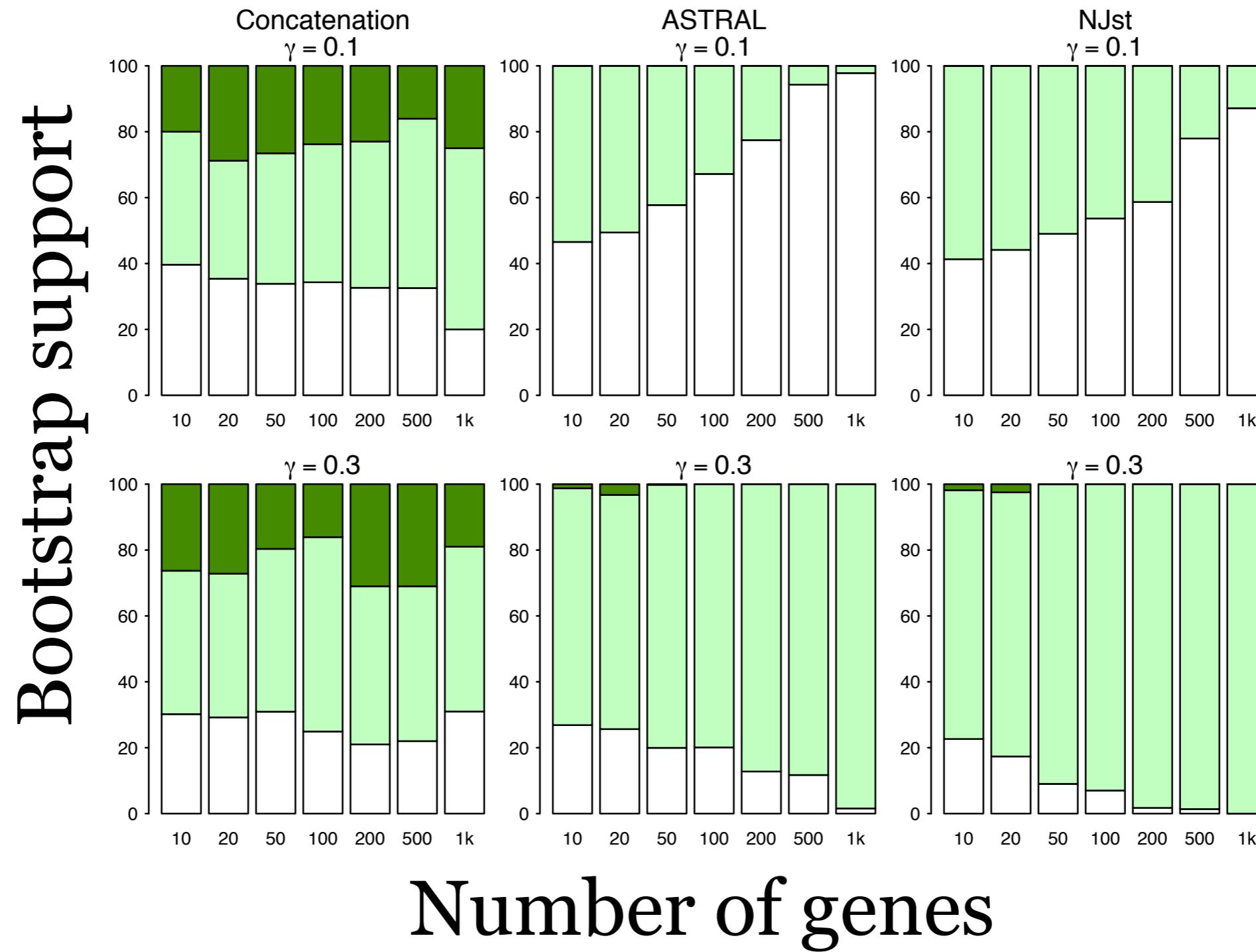
# Why?

Phylogenetic network

Coalescent tree methods  
not robust to gene flow



White:  
true tree



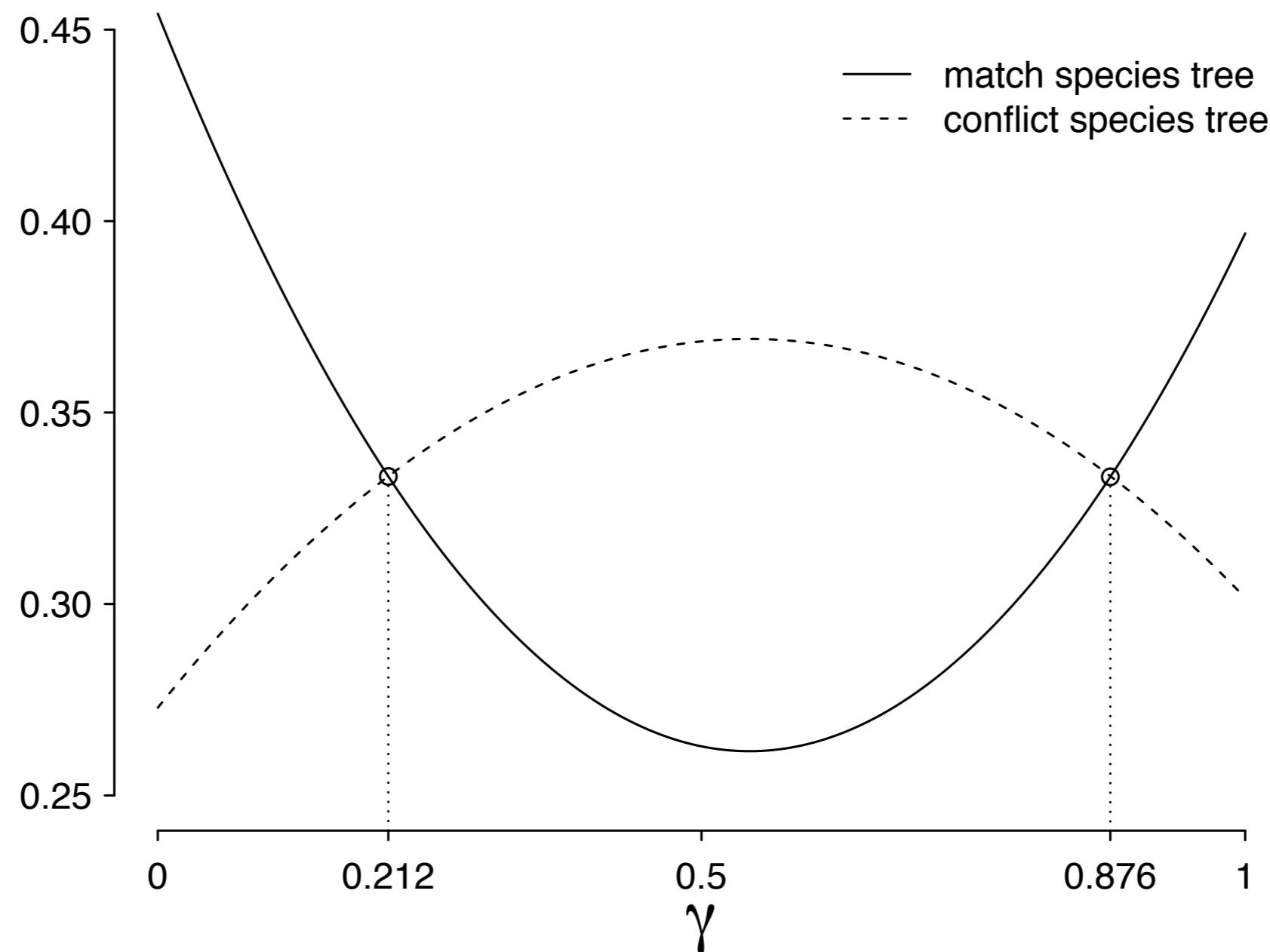
Number of genes

(S.-L., Yang, Ané, 2016, Syst Bio)

ASTRAL (Mirarab et al, 2014)  
NJst (Liu&Yu, 2011)

# Why? Phylogenetic network

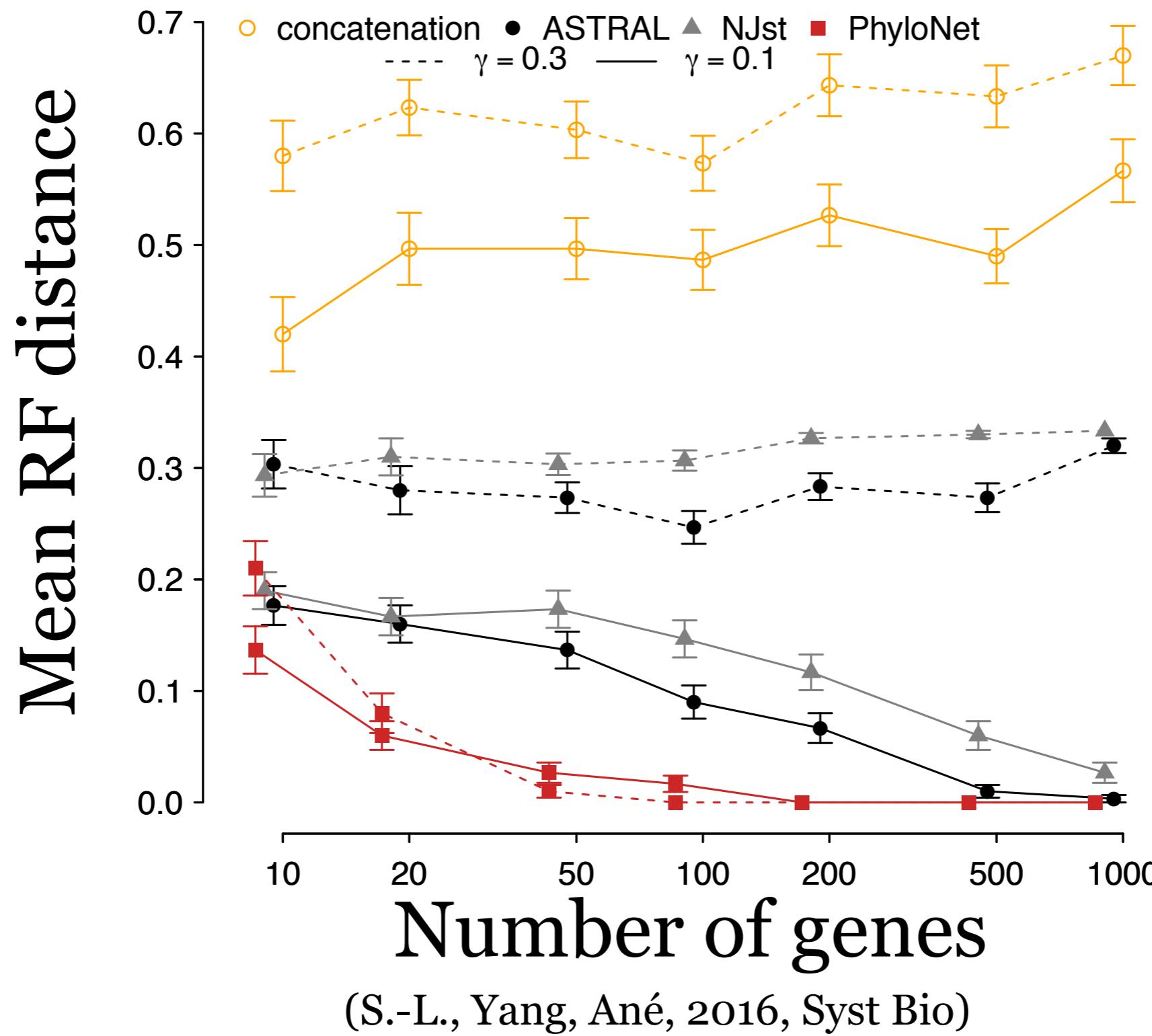
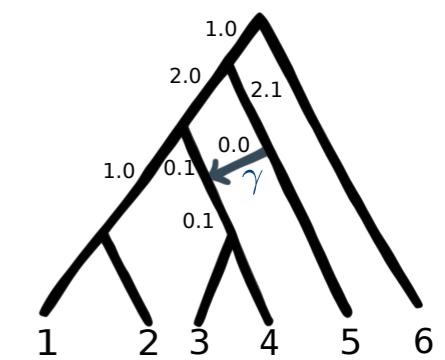
## Anomaly zone with gene flow



(S.-L., Yang, Ané, 2016, Syst Bio)

# Why? Phylogenetic network

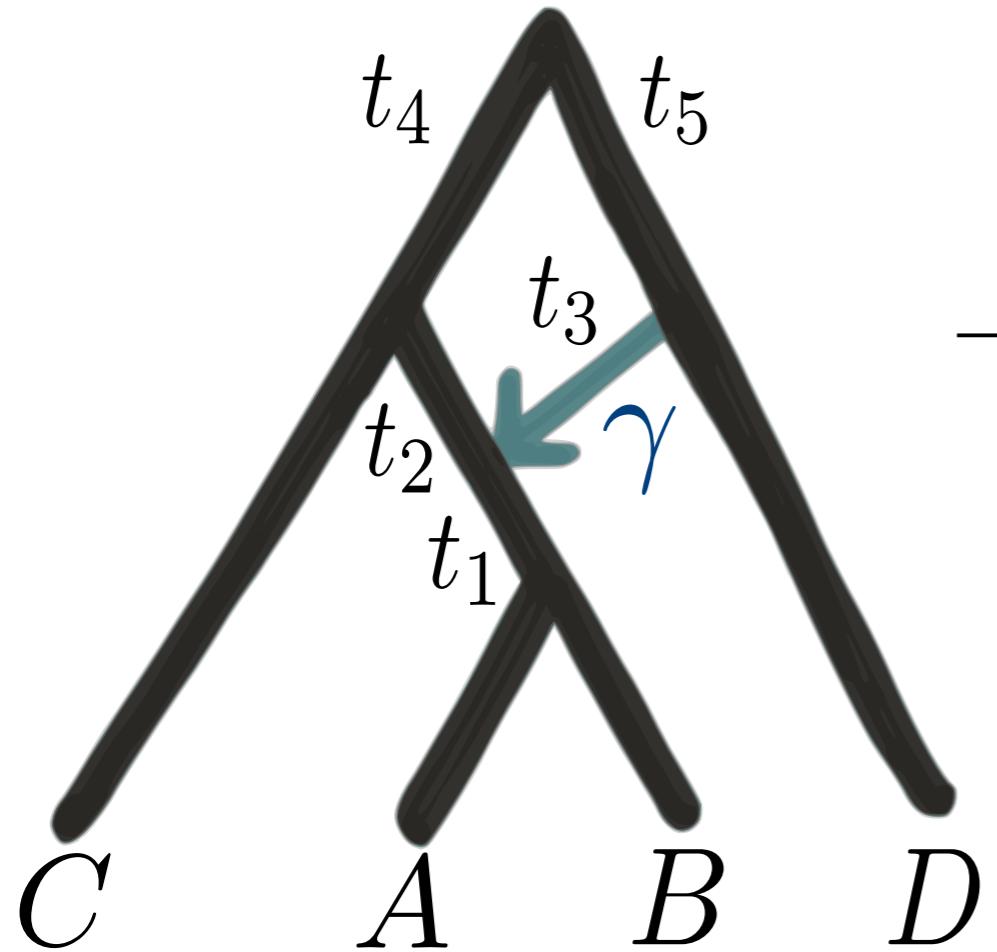
Coalescent tree methods  
not robust to gene flow



# Why?

Phylogenetic network

Anomalous unrooted  
gene trees with gene flow



Frequency among gene trees

Quartet	$\gamma = 0.0$	$\gamma = 0.1$	$\gamma = 0.3$
$AB CD$	<b>0.347</b>	0.298	0.260
$CA BD$	0.327	0.351	0.370
$CB AD$	0.327	0.351	0.370

$$t_1 = t_2 = 0.01, t_3 = t_4 = t_5 = 1$$

- **ILS**: no AUGT on 4 taxa (Degnan, 2013)
- **ILS+HGT**: AUGT on 4 taxa (S.-L., Yang, Ané, 2016, Syst Bio)

# So far...

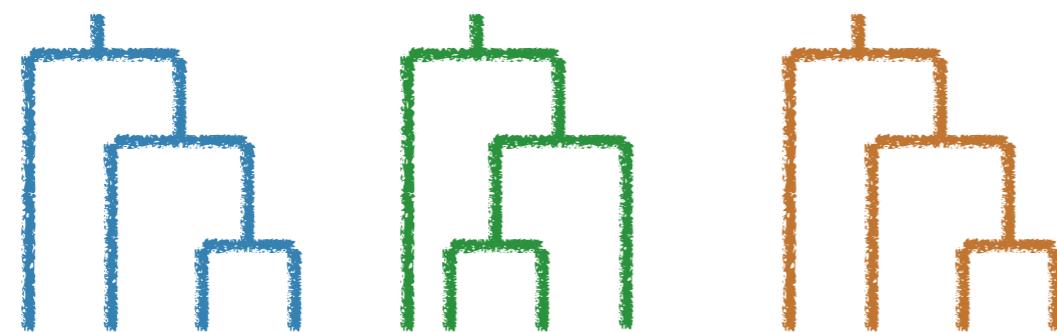
- Networks are good
- Explicit networks are better
- If you ignore gene flow, you can estimate the wrong tree

# How?

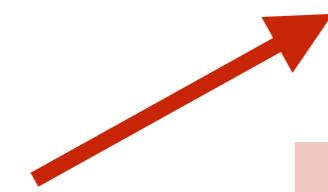
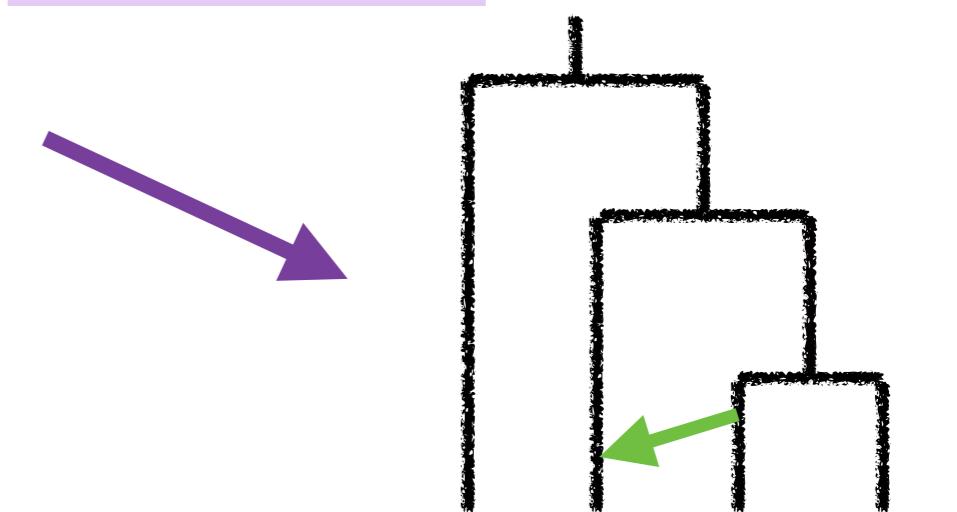
## Phylogenetic network



MrBayes  
(Huelsenbeck, Ronquist, 2001)  
RAxML  
(Stamatakis, 2014)  
PhyML  
(Guindon et al, 2010)



BEAST2  
(Zhang et al, 2017)  
PhyloNet  
(Wen et al, 2016)



SNaQ  
(S.-L., Ane, 2016)  
PhyloNet  
(Yu et al, 2014)

# Multispecies coalescent on a network



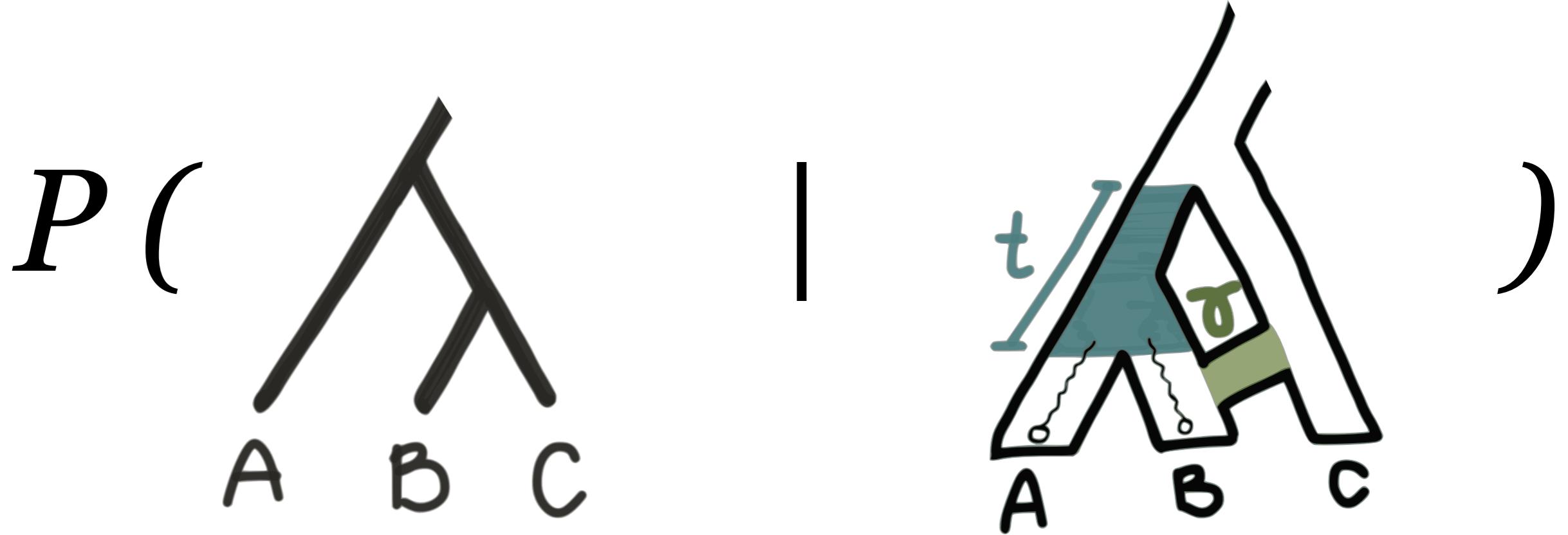
(Meng, Kubatko, 2009)  
(Yu, Degnan, Nakhleh, 2012)

# Multispecies coalescent on a network



(Meng, Kubatko, 2009)  
(Yu, Degnan, Nakhleh, 2012)

# Multispecies coalescent on a network



(Meng, Kubatko, 2009)  
(Yu, Degnan, Nakhleh, 2012)



<https://solislemuslab.github.io/>

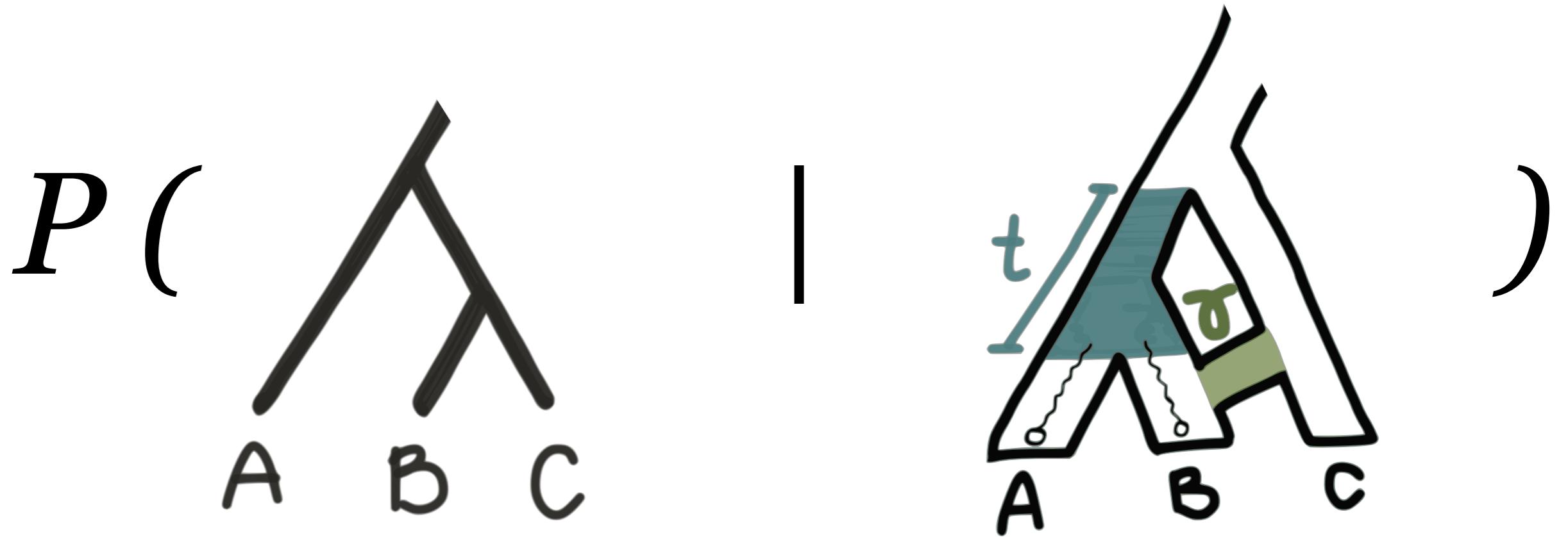


@solislemuslab



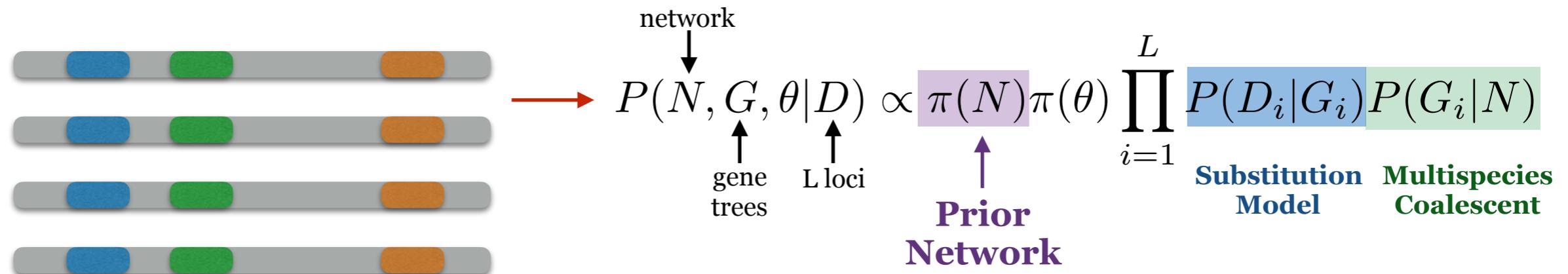
crsl4

# Multispecies coalescent on a network



$$p_{BC|AD}(t, t_2, \gamma) = (1 - \gamma) \frac{1}{3} e^{-t} + \gamma (1 - \frac{2}{3} e^{-t_2})$$

(Meng, Kubatko, 2009)  
(Yu, Degnan, Nakhleh, 2012)

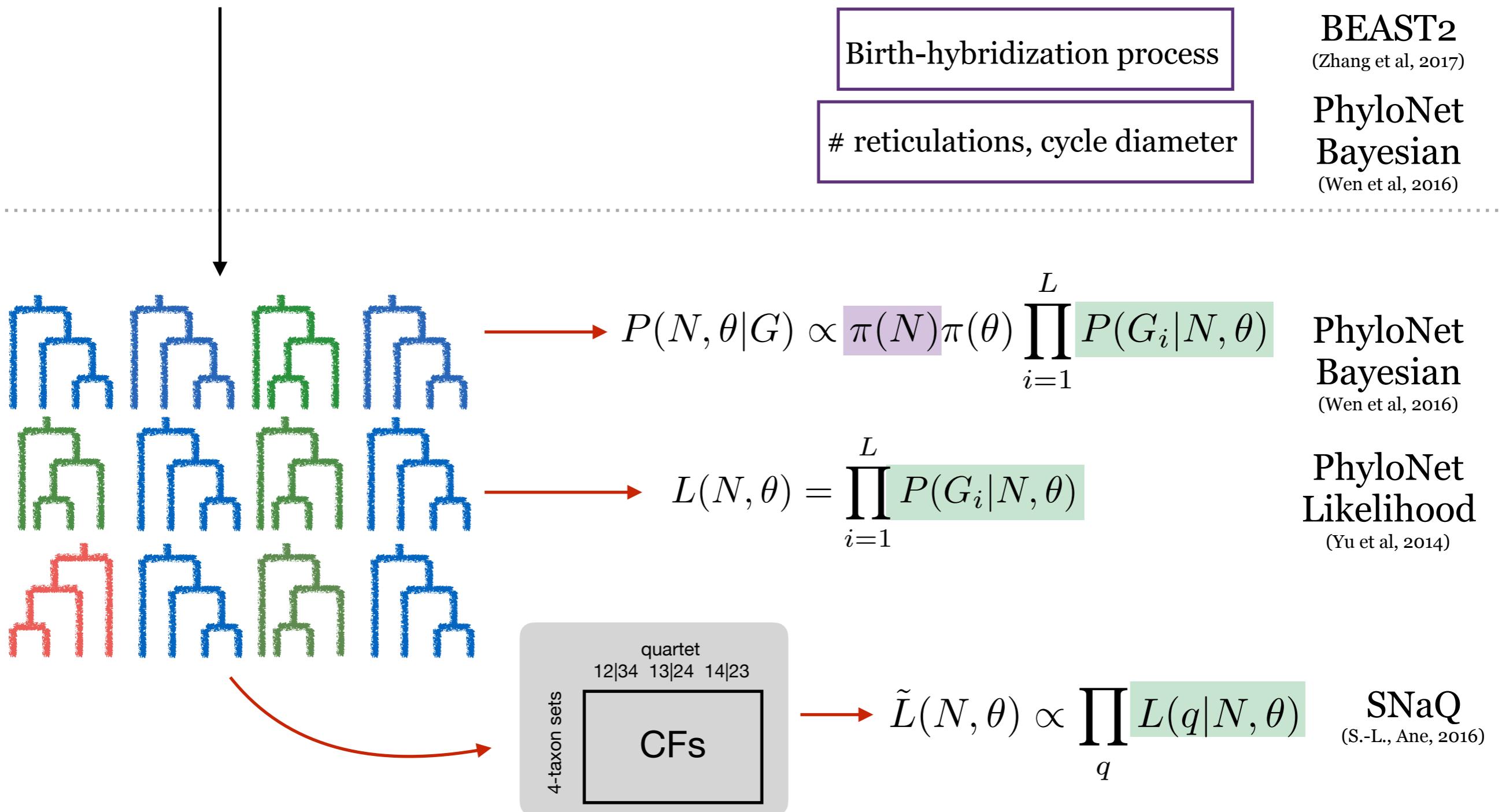
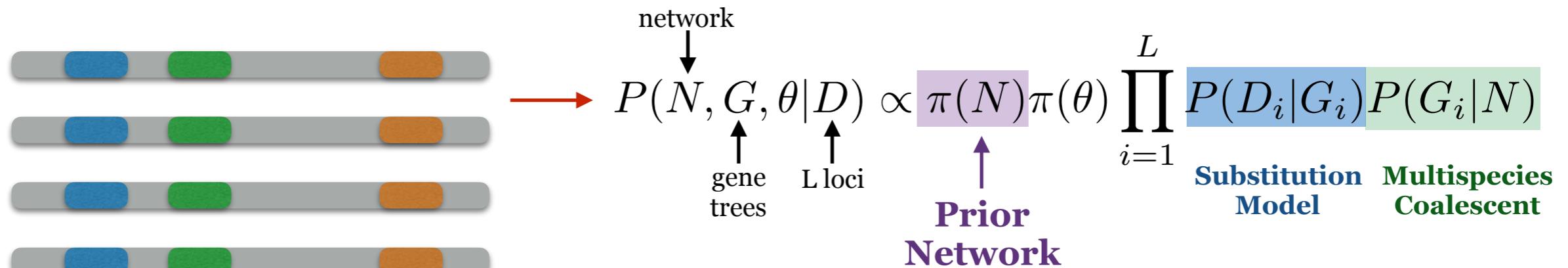


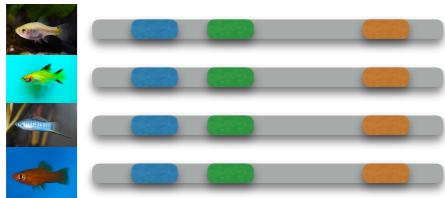
Birth-hybridization process

# reticulations, cycle diameter

**BEAST2**  
(Zhang et al, 2017)

**PhyloNet Bayesian**  
(Wen et al, 2016)

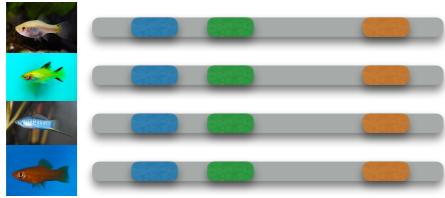




BEAST2  
(Zhang et al, 2017)

Birth-hybridization process

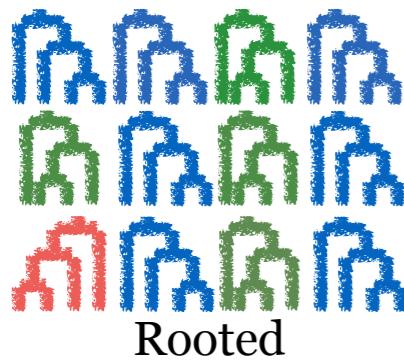
Most accurate,  
not scalable



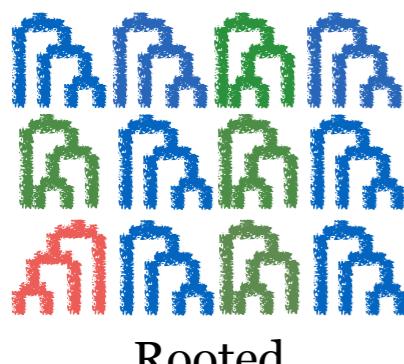
PhyloNet  
Bayesian  
(Wen et al, 2016)

**MCMC:**  
Network moves,  
mixing

# reticulations,  
cycle diameter

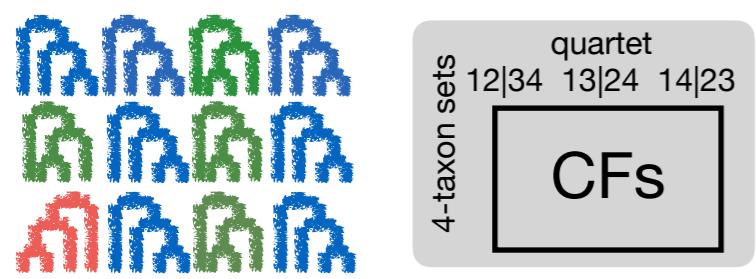


PhyloNet  
Bayesian  
(Wen et al, 2016)



PhyloNet  
Likelihood  
(Yu et al, 2014)

**Heuristic search:**  
Network moves



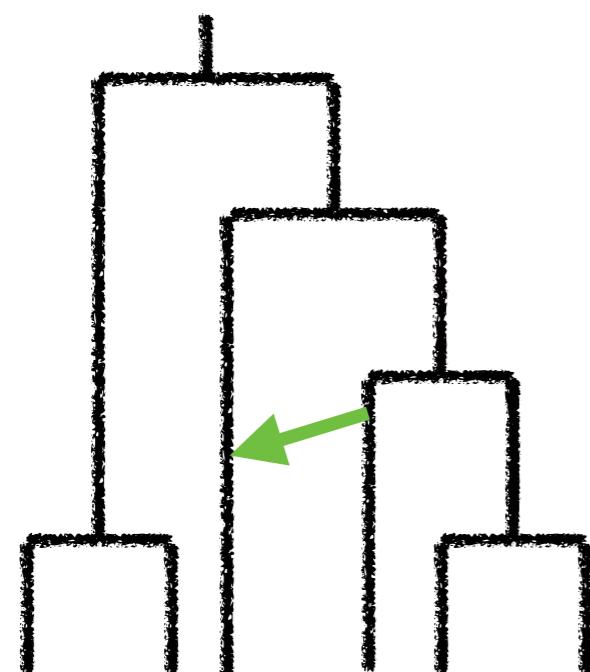
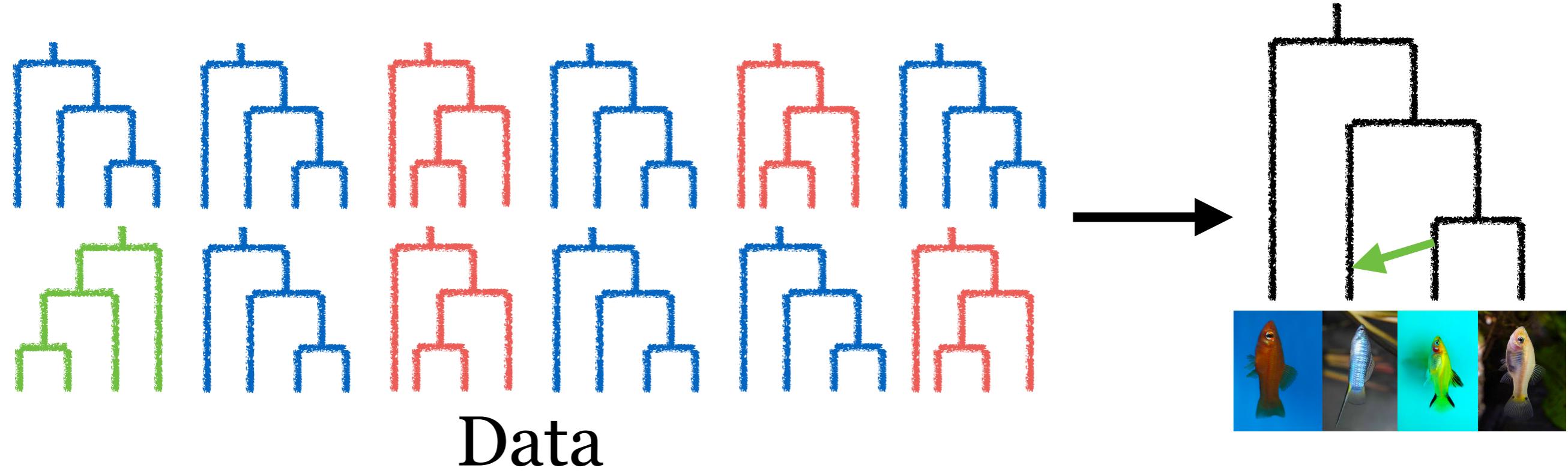
SNaQ  
(S.-L., Ane, 2016)

Level-1  
networks

More scalable,  
Robust

STEM-hy	gene trees rooted, BL	likelihood	hybridization b/w sister lineages
PhyloNet InferNetwork_ML	gene trees rooted	likelihood	
PhyloNet InferNetwork_MPL	gene trees rooted	triplet likelihood	
Phylogenetworks SNaQ	gene trees or quartet CFs	quartet likelihood	level-1 network
PhyloNet MCMC_GT	gene trees rooted	Bayesian	compound prior
PhyloNet MCMC_SEQ	alignments	Bayesian	compound prior no rate variation
BEAST2 SpeciesNetwork	alignments	Bayesian	birth-hyb prior
PhyloNet MLE_BiMarkers	biallelic sites	likelihood	compound prior
PhyloNet MCMC_BiMarkers	biallelic sites	Bayesian	compound prior
HyDe	sites	invariants	4 taxa, 1 hyb.

# Maximum pseudolikelihood



$$\tilde{L}(\textit{network}) = \prod L(\textit{quartet})$$

(S-L, Ané, 2016, PLoS Genetics)

[www.github.com/CRSL4/PhyloNetworks](https://www.github.com/CRSL4/PhyloNetworks)

Quartet-based inference

snaQ julia



<https://solislemuslab.github.io/>

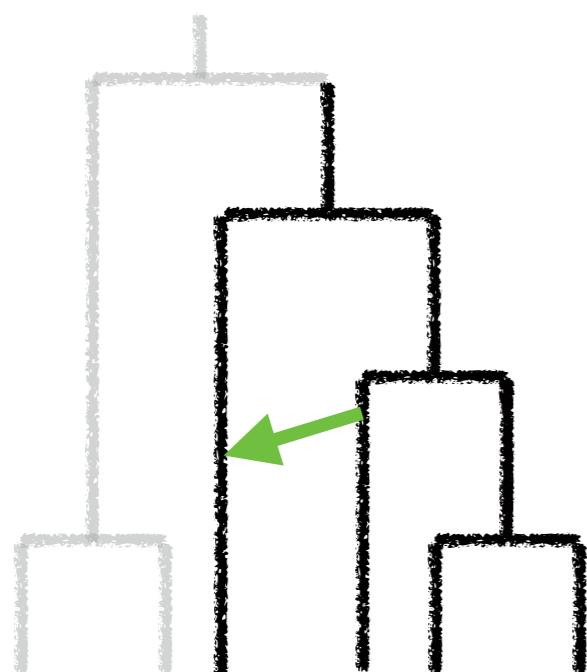
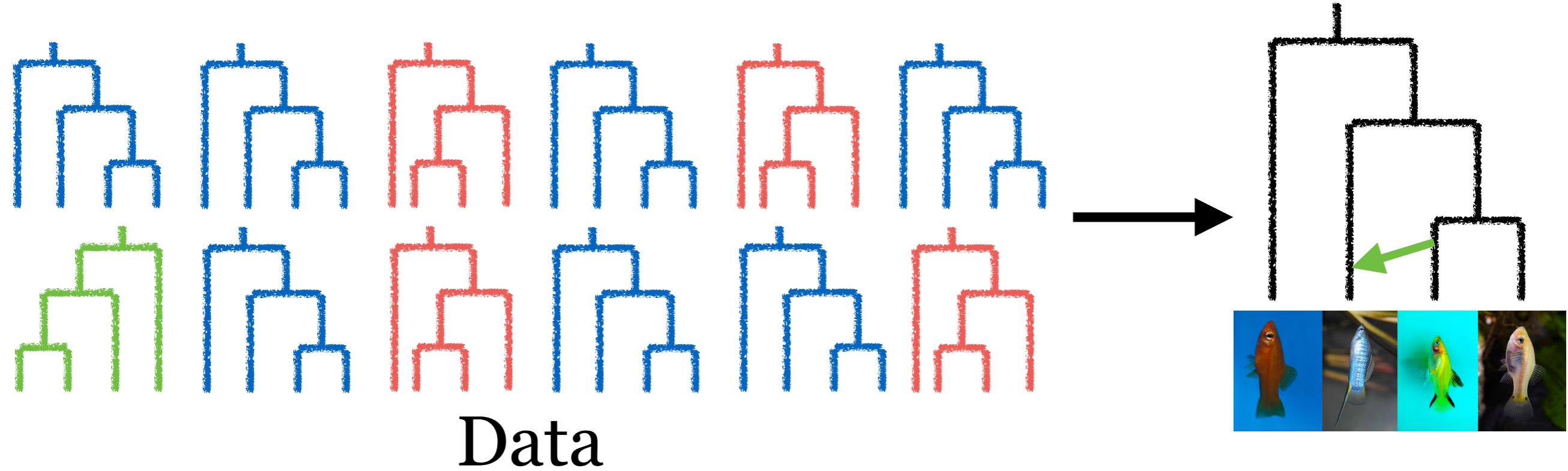


@solislemuslab



crsl4

# Maximum pseudolikelihood



Quartet-based inference

$$\tilde{L}(\textit{network}) = \prod L(\textit{quartet})$$

(S-L, Ané, 2016, PLoS Genetics)

[www.github.com/CRSL4/PhyloNetworks](https://github.com/CRSL4/PhyloNetworks)

snaQ julia



<https://solislemuslab.github.io/>

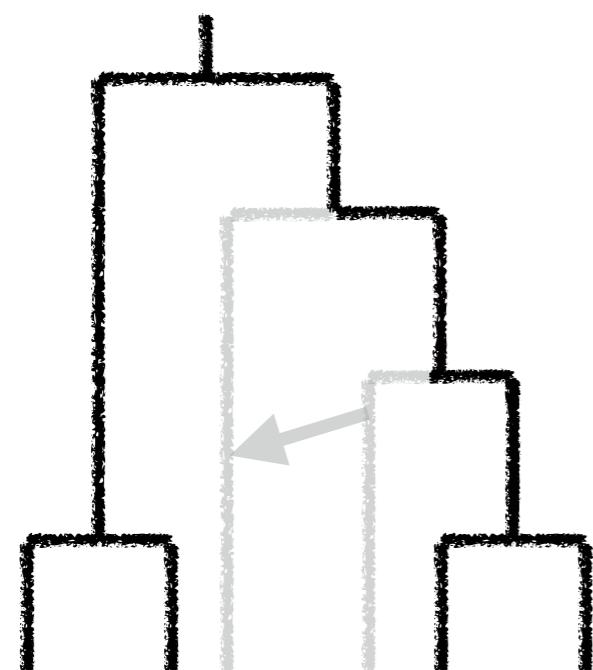
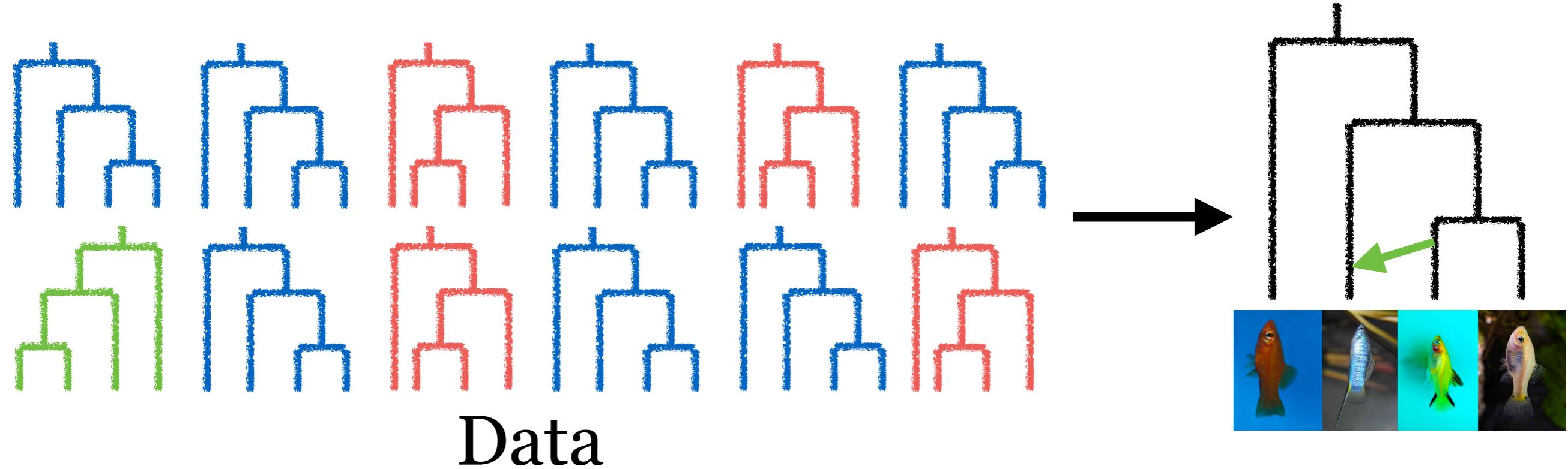


@solislemuslab



crsl4

# Maximum pseudolikelihood



Quartet-based inference

$$\tilde{L}(\textit{network}) = \prod L(\textit{quartet})$$

(S-L, Ané, 2016, PLoS Genetics)

[www.github.com/CRSL4/PhyloNetworks](https://github.com/CRSL4/PhyloNetworks)

snaQ julia



<https://solislemuslab.github.io/>



@solislemuslab



crsl4

# Maximum pseudolikelihood

Unrooted gene trees

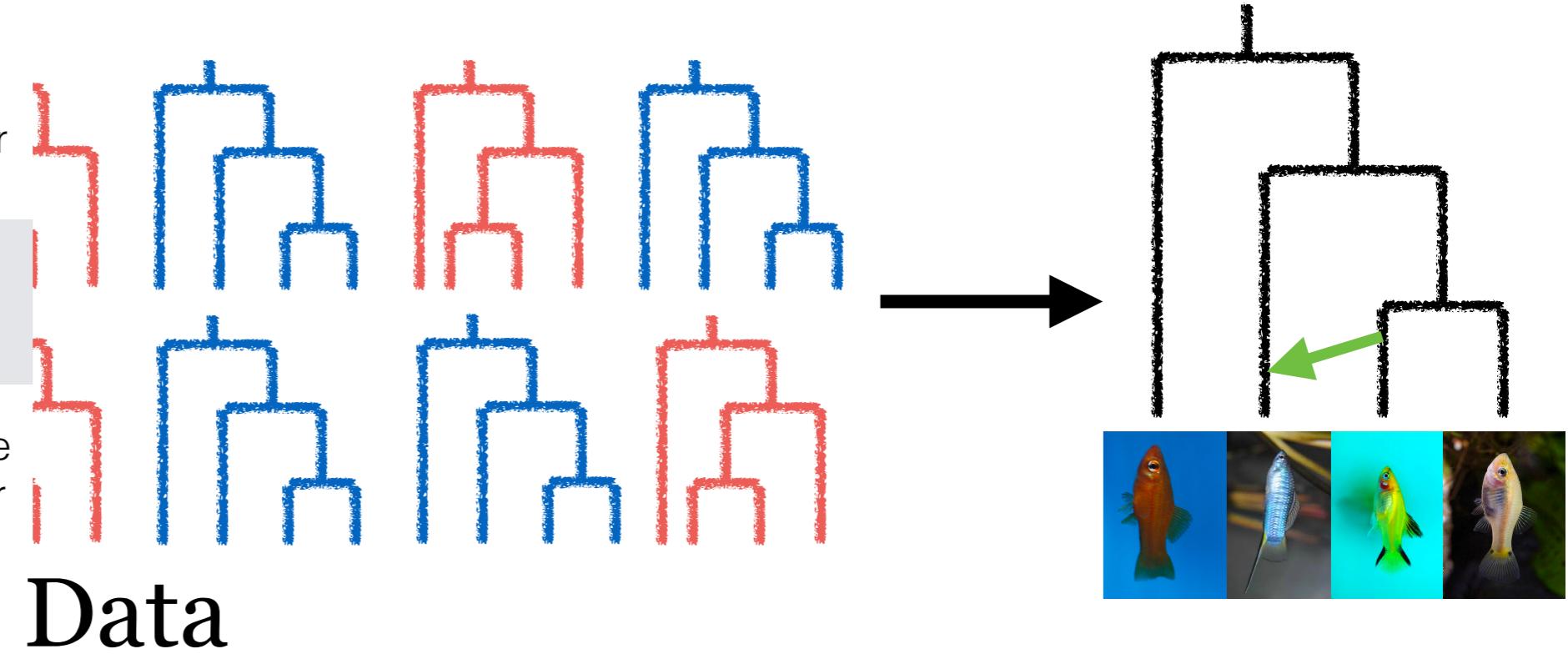
No branch lengths

Concordance factors

No rooting error

No molecular clock assumption

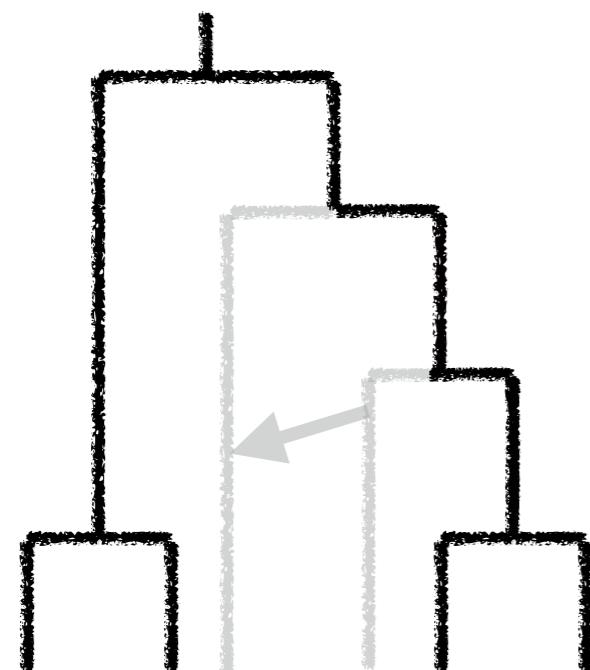
Account for tree estimation error



$$\tilde{L}(\textit{network}) = \prod L(\textit{quartet})$$

(S-L, Ané, 2016, PLoS Genetics)

[www.github.com/CRSL4/PhyloNetworks](https://github.com/CRSL4/PhyloNetworks)



Quartet-based inference

snaQ julia



<https://solislemuslab.github.io/>

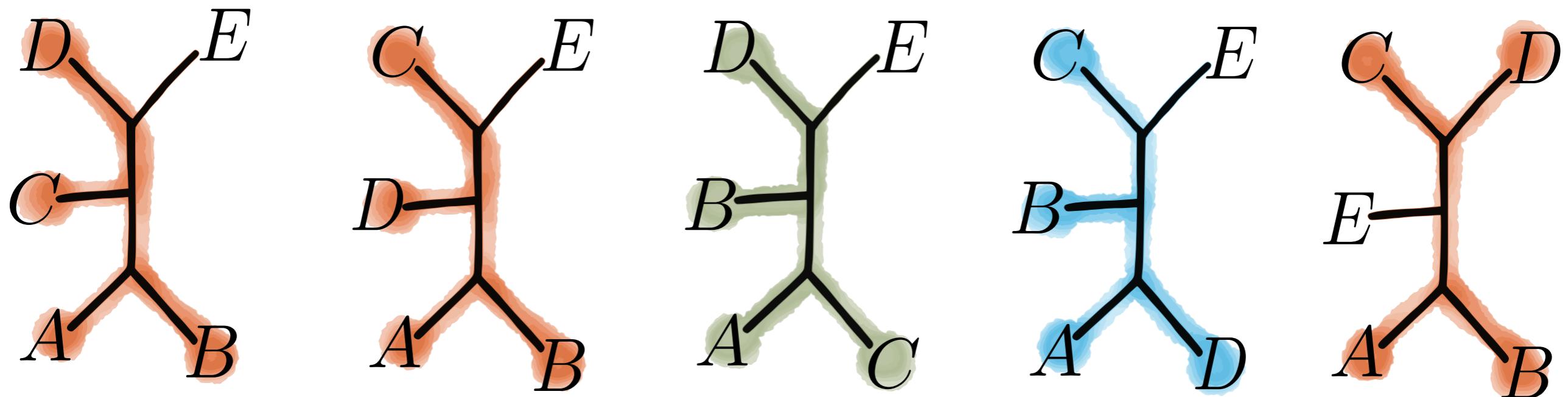


@solislemuslab

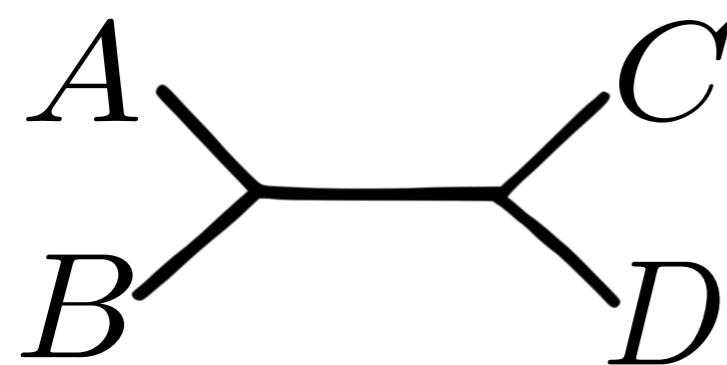


crsl4

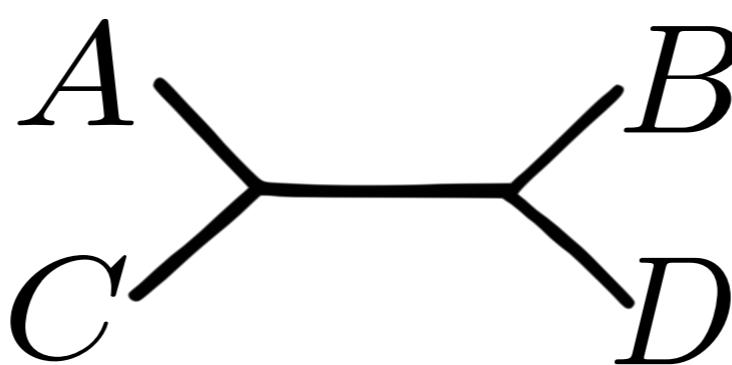
# Quartet-based inference



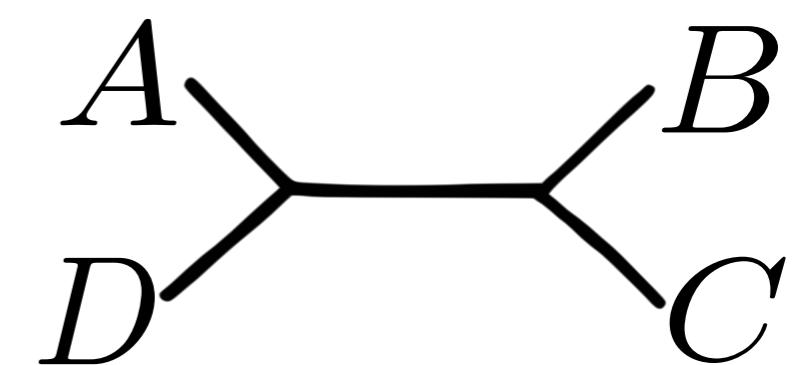
Concordance factors (CF):  
% of genes having the quartet in their tree



3/5



1/5



1/5



<https://solislemuslab.github.io/>



@solislemuslab



crsl4

# Quartet-based inference

Observed **quartet** CFs:

4 taxon set	$CF_1$	$CF_2$	$CF_3$
A B C D	.80	.10	.10
A B C E	.40	.40	.20
A B D E	.40	.40	.20
A C D E	.84	.08	.08
B C D E	.82	.10	.08

inferred network:

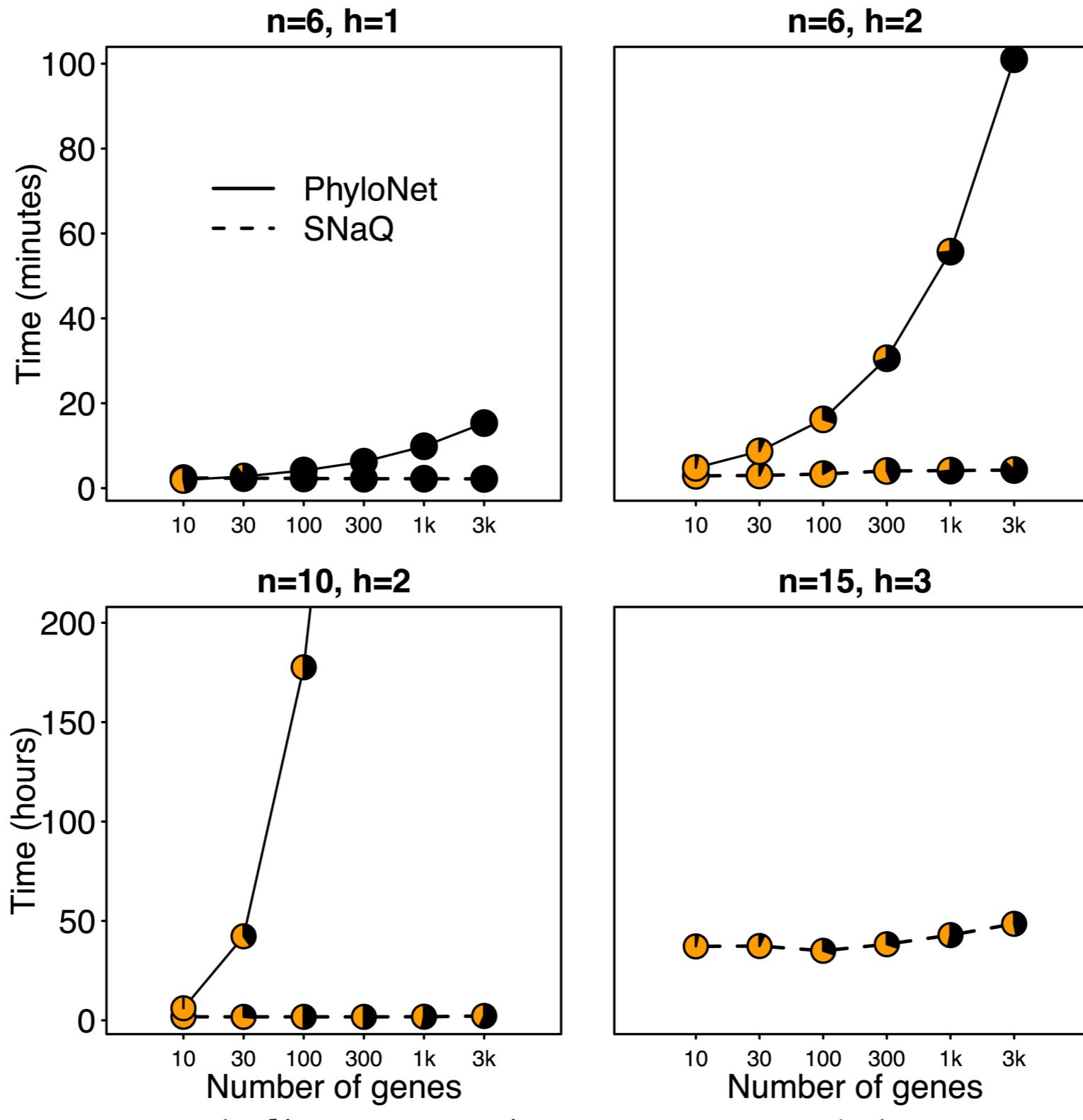


Maximum Pseudo-Likelihood:

$$\log \tilde{L} = \sum_{q \in Q(N)} CF_{in,1} \log(CF_{net,1}) + CF_{in,2} \log(CF_{net,2}) + CF_{in,3} \log(CF_{net,3})$$



# Scalability gains



<https://solislemuslab.github.io/>



@solislemuslab

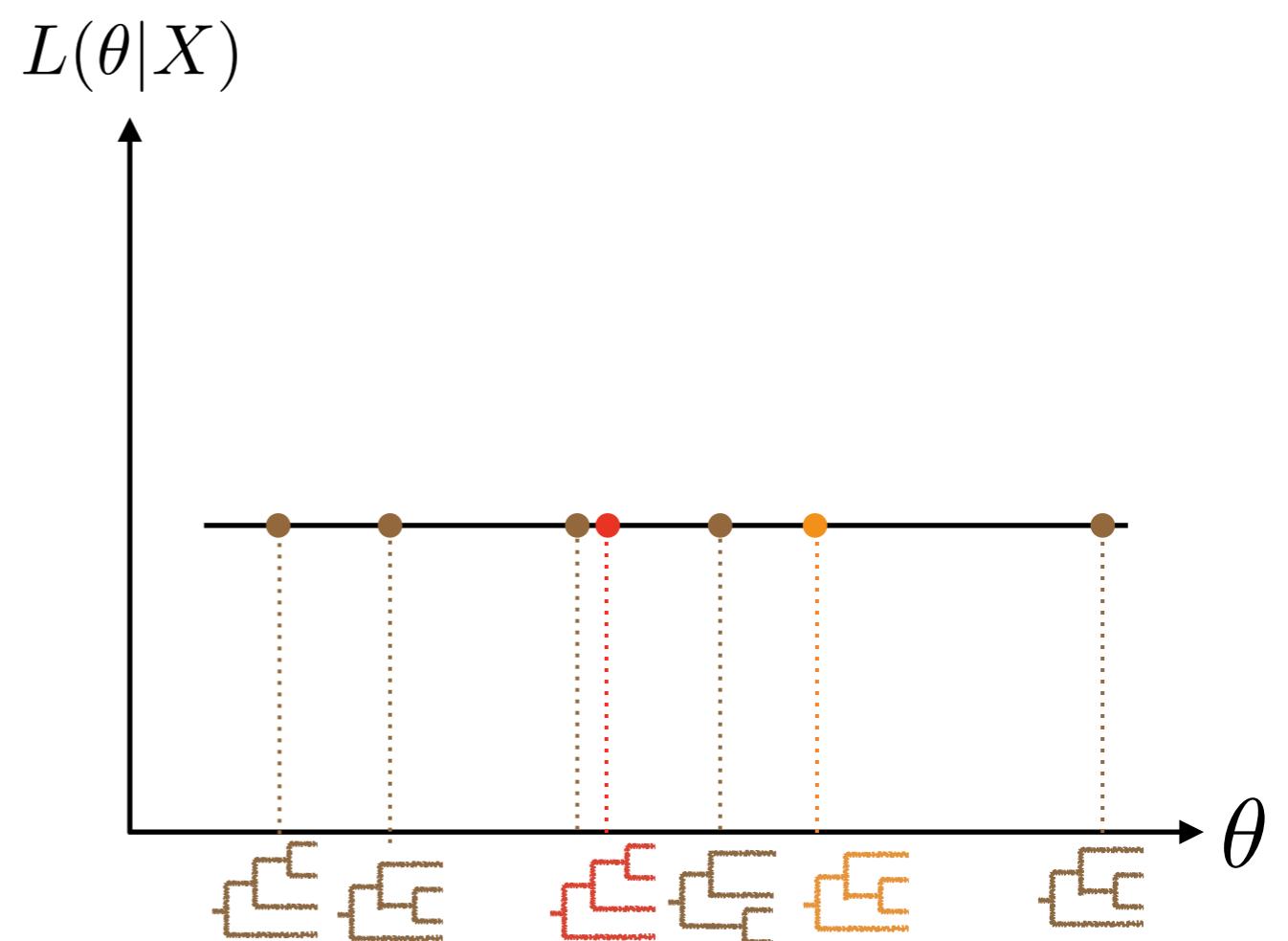
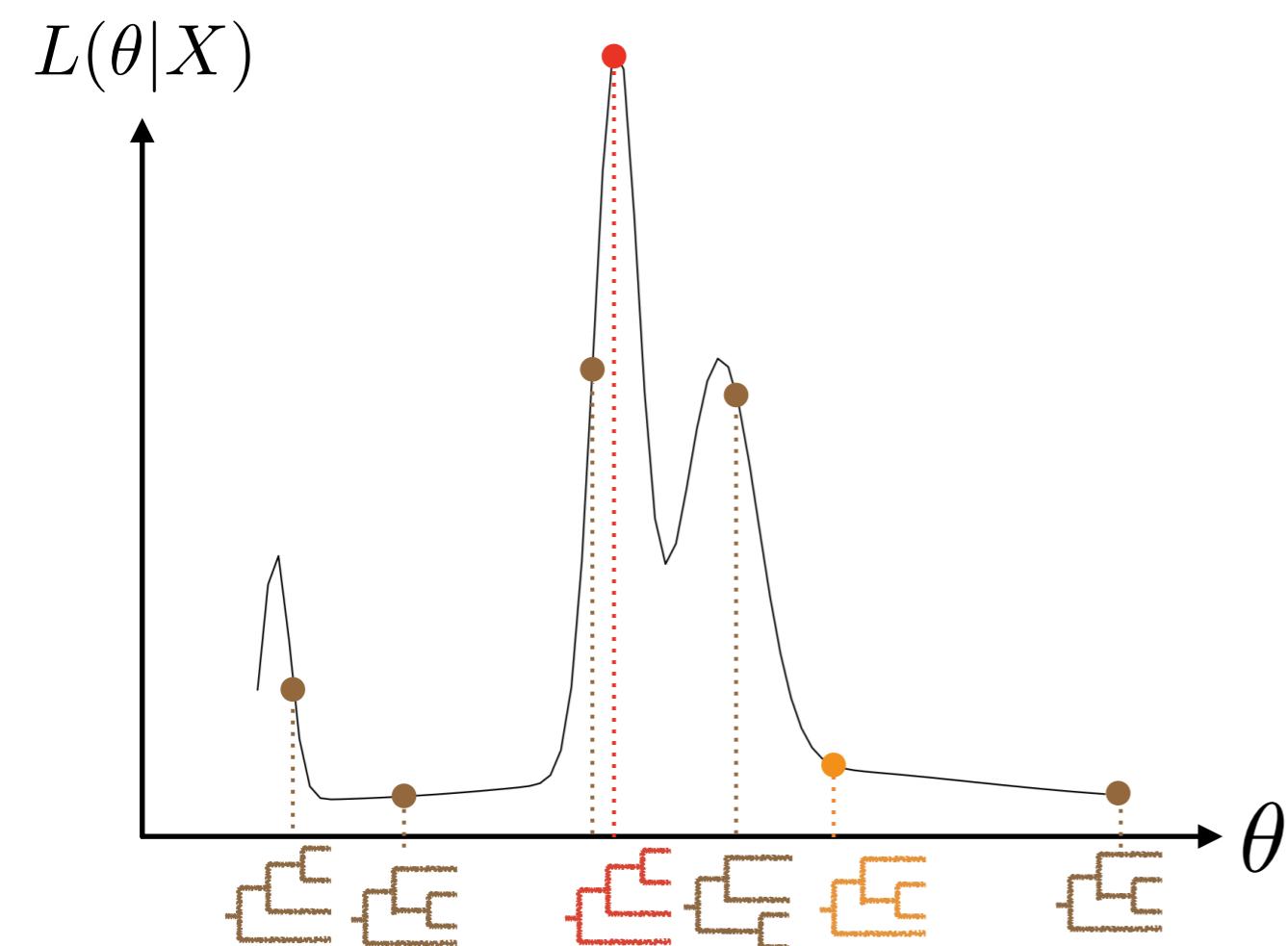


crsl4

# Challenges

- Network space
- Identifiability
- Network comparison

# Identifiability

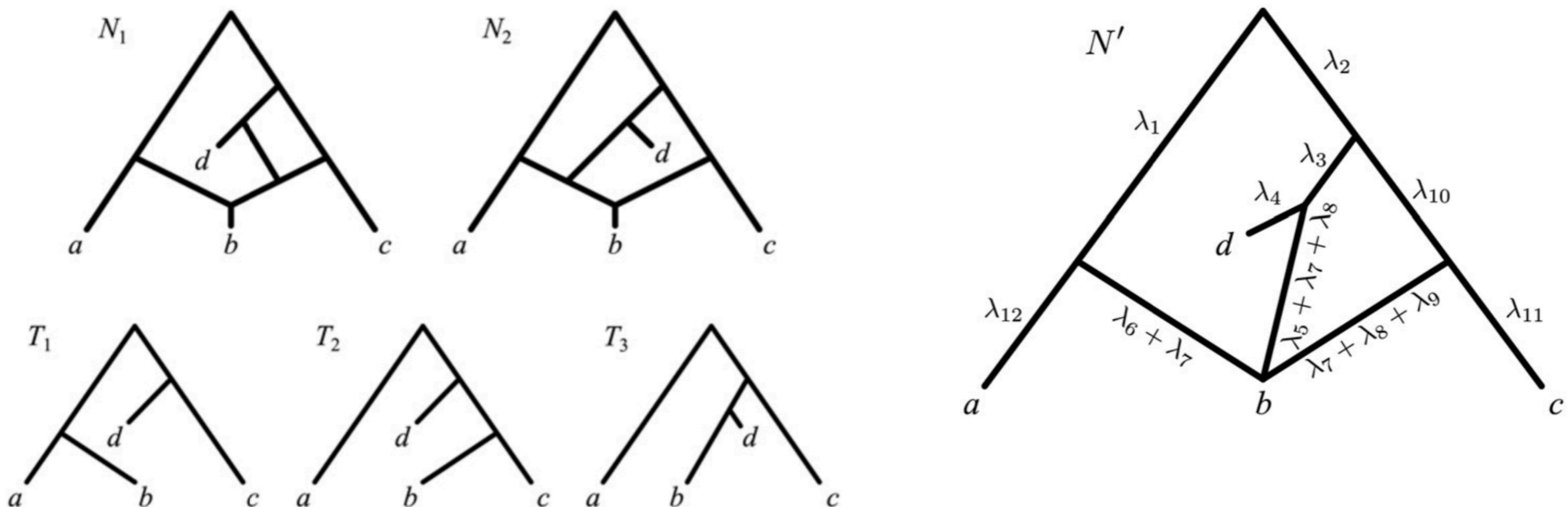


RESEARCH ARTICLE

# Reconstructible Phylogenetic Networks: Do Not Distinguish the Indistinguishable

Fabio Pardi<sup>1,3\*</sup>, Celine Scornavacca<sup>2,3</sup>

**1** Laboratoire d’Informatique, de Robotique et de Microélectronique de Montpellier (LIRMM, UMR 5506) CNRS, Université de Montpellier, France, **2** Institut des Sciences de l’Evolution de Montpellier (ISE-M, UMR 5554) CNRS, IRD, Université de Montpellier, France, **3** Institut de Biologie Computationnelle, Montpellier, France

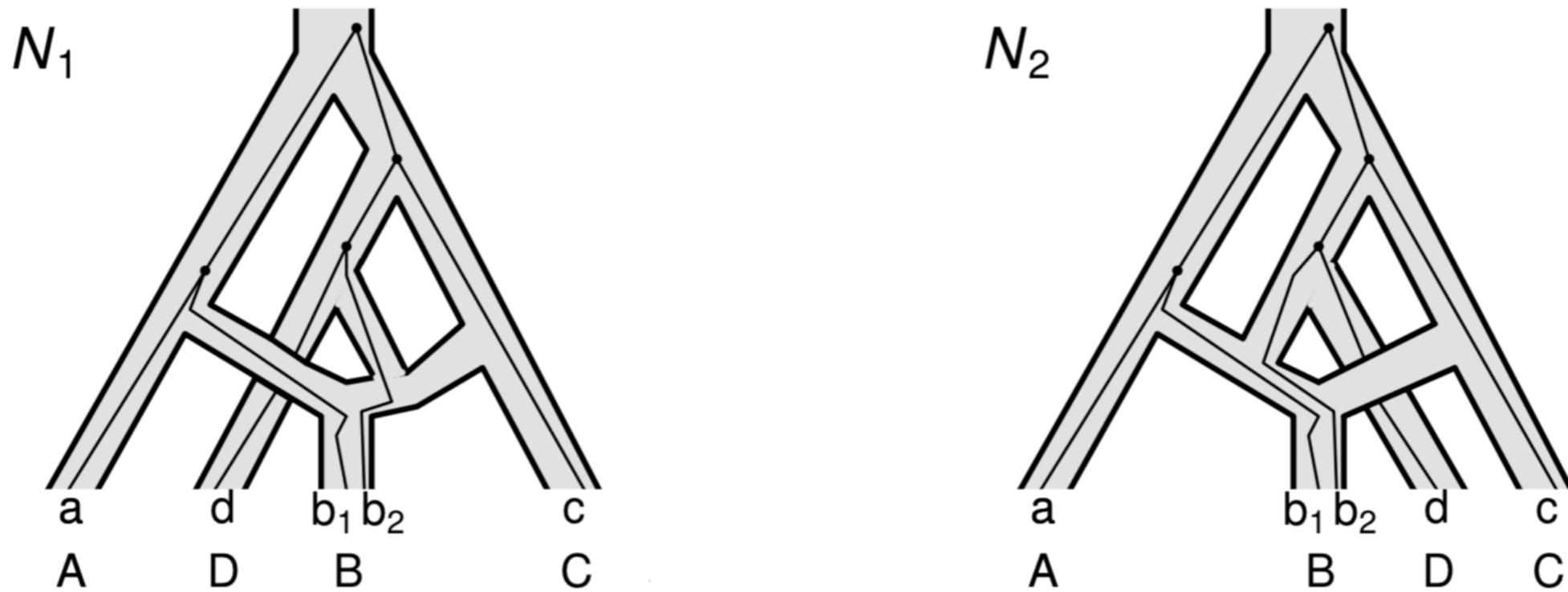


Undistinguishable with the  
“displayed trees” criterion

Solution: Canonical  
network (“unzipped”)

# Displayed Trees Do Not Determine Distinguishability Under the Network Multispecies Coalescent

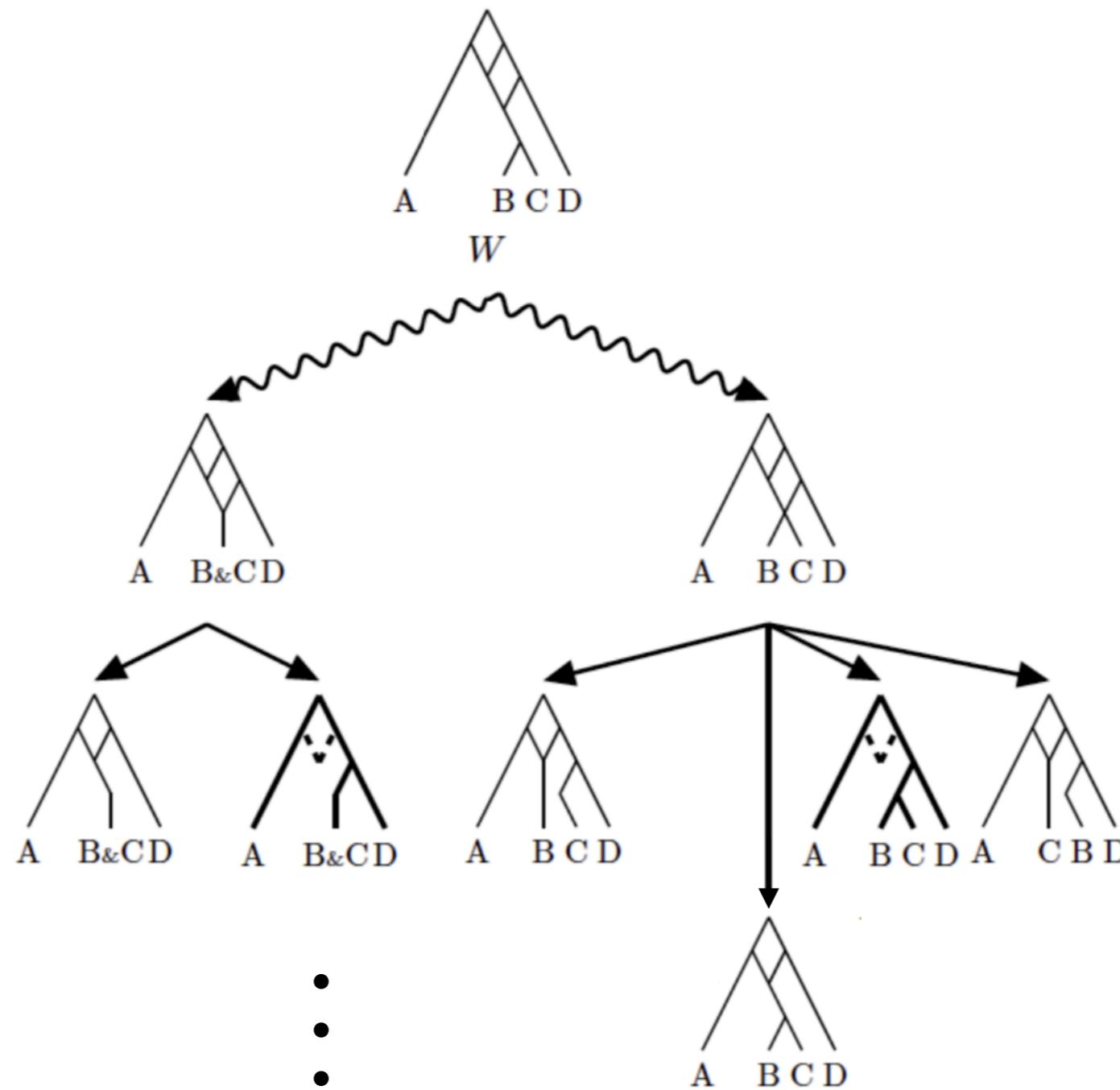
Sha Zhu<sup>1</sup>, James H. Degnan<sup>2</sup>



Distinguishable under the MSC

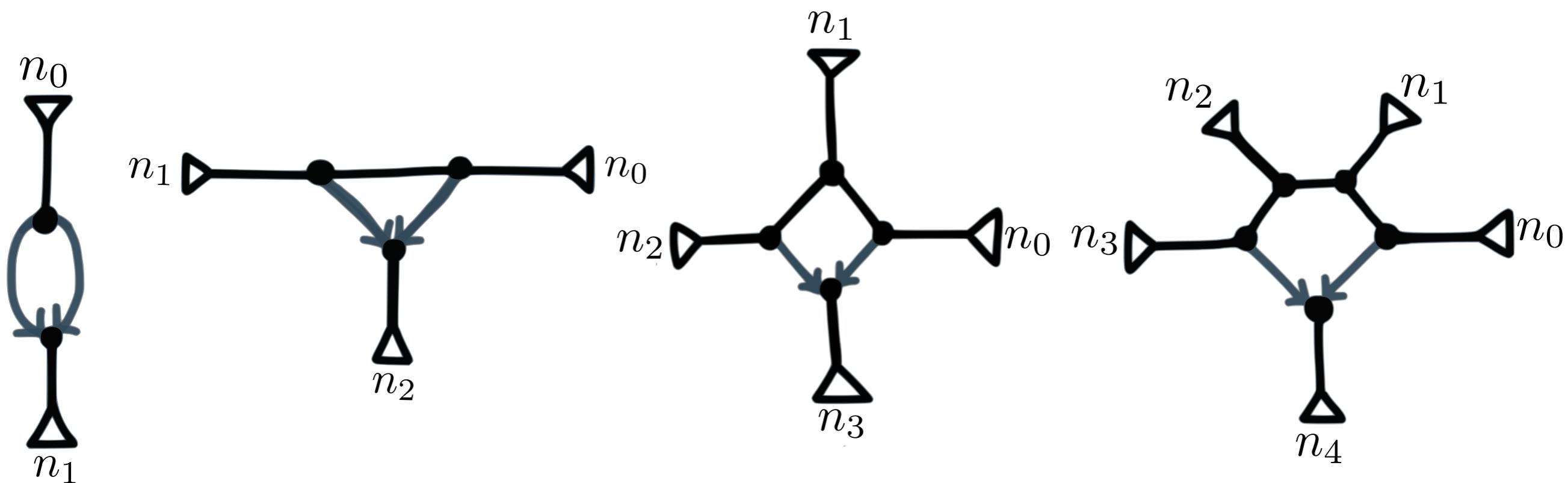
# Displayed Trees Do Not Determine Distinguishability Under the Network Multispecies Coalescent

Sha Zhu<sup>1</sup>, James H. Degnan<sup>2</sup>



Decomposing network in **parental** trees

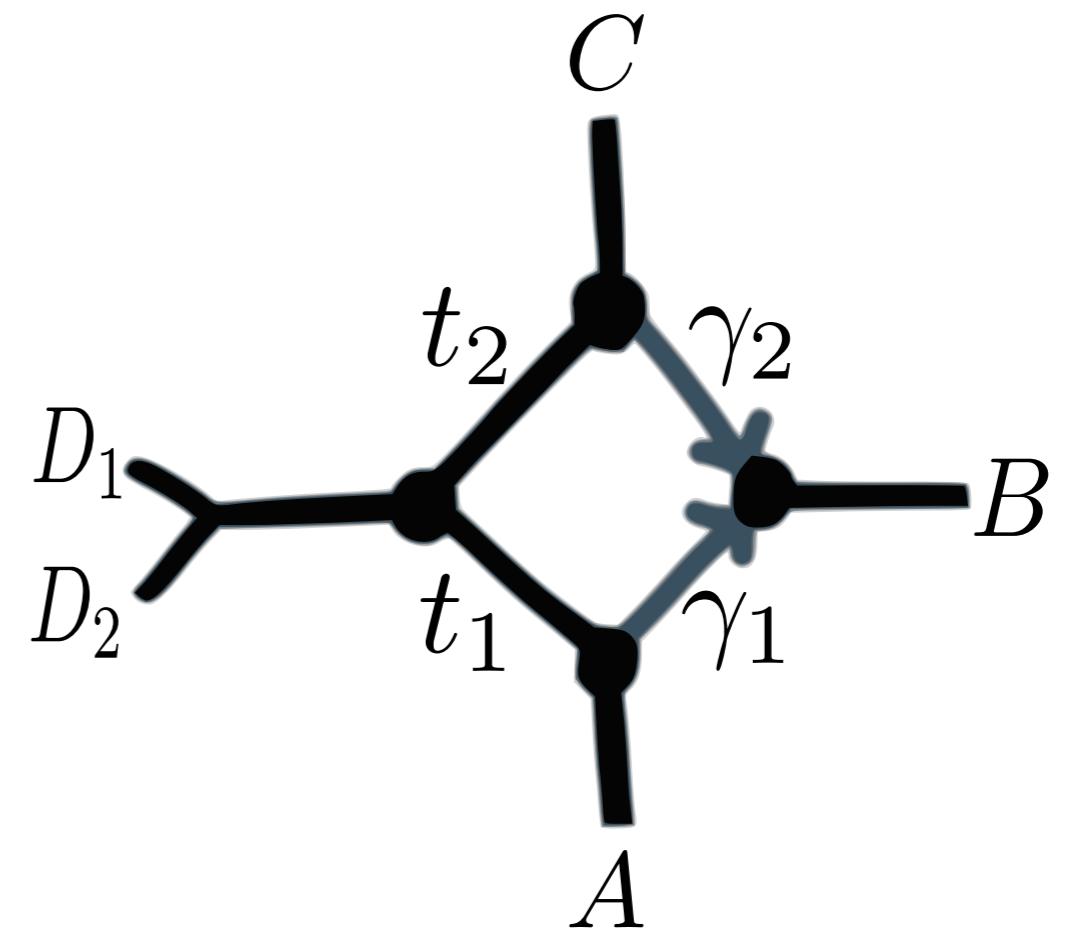
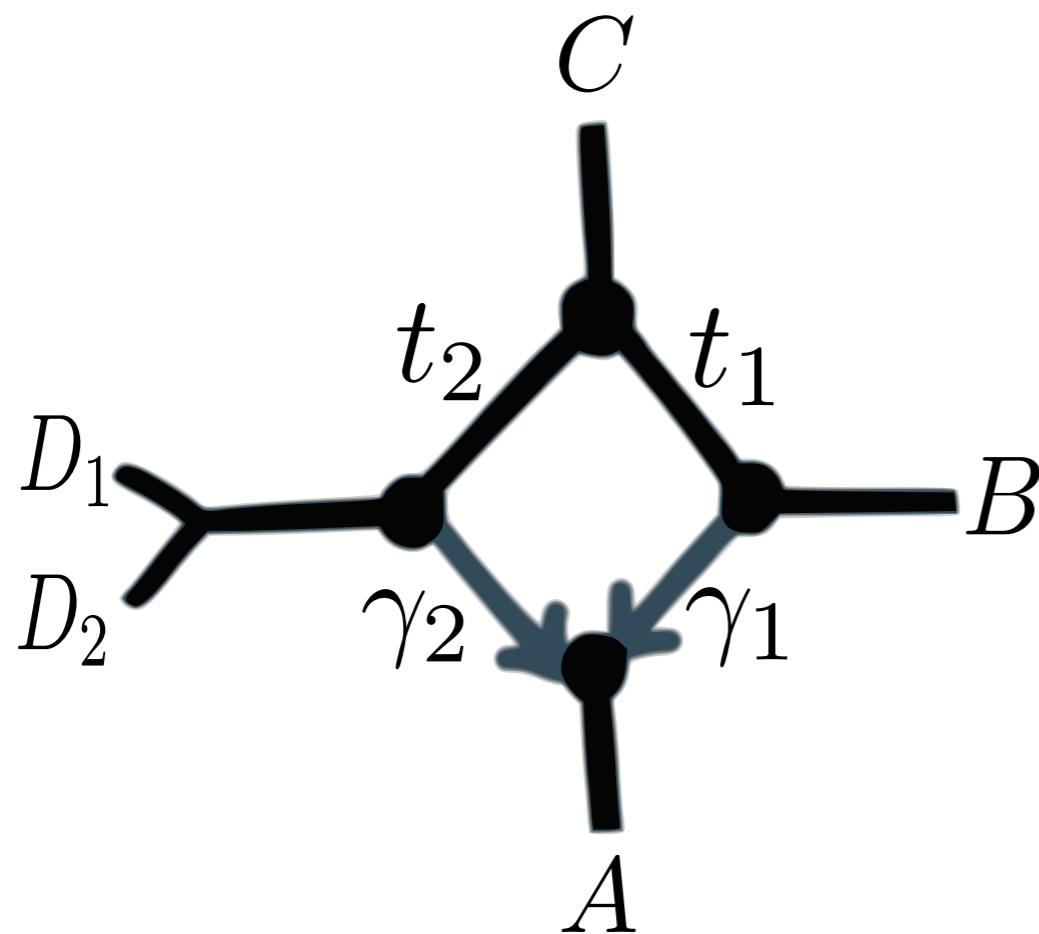
## RESEARCH ARTICLE

Inferring Phylogenetic Networks with  
Maximum Pseudolikelihood under  
Incomplete Lineage SortingClaudia Solís-Lemus<sup>1\*</sup>, Cécile Ané<sup>1,2</sup>Can we detect the  
presence of  
hybridization in level-1  
networks?**No****Yes**  
 $(n_i, n_j \geq 2)$ **Yes**  
 $(n_i \geq 2)$ **Yes**

Generic Identifiability

 $t_i \in (0, \infty), \gamma \in (0, 1)$

# In practice: flat pseudolikelihood

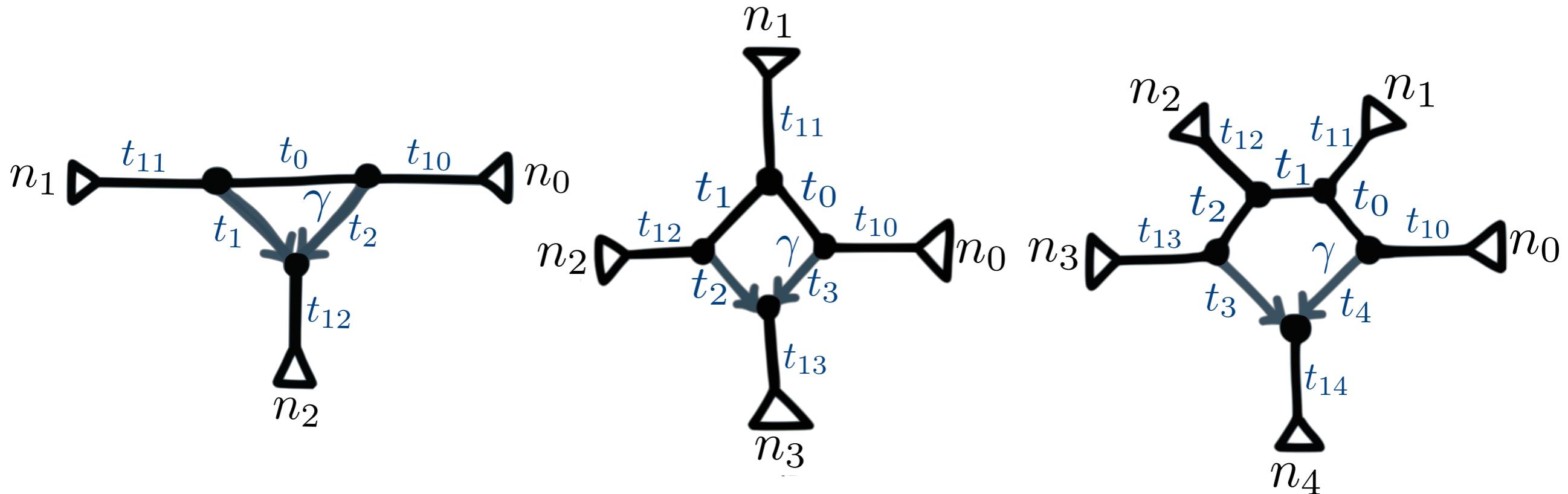


# Can we estimate numerical parameters?

RESEARCH ARTICLE

## Inferring Phylogenetic Networks with Maximum Pseudolikelihood under Incomplete Lineage Sorting

Claudia Solís-Lemus<sup>1\*</sup>, Cécile Ané<sup>1,2</sup>



No

Good triangle  
( $t_{12} = 0$ )

Yes

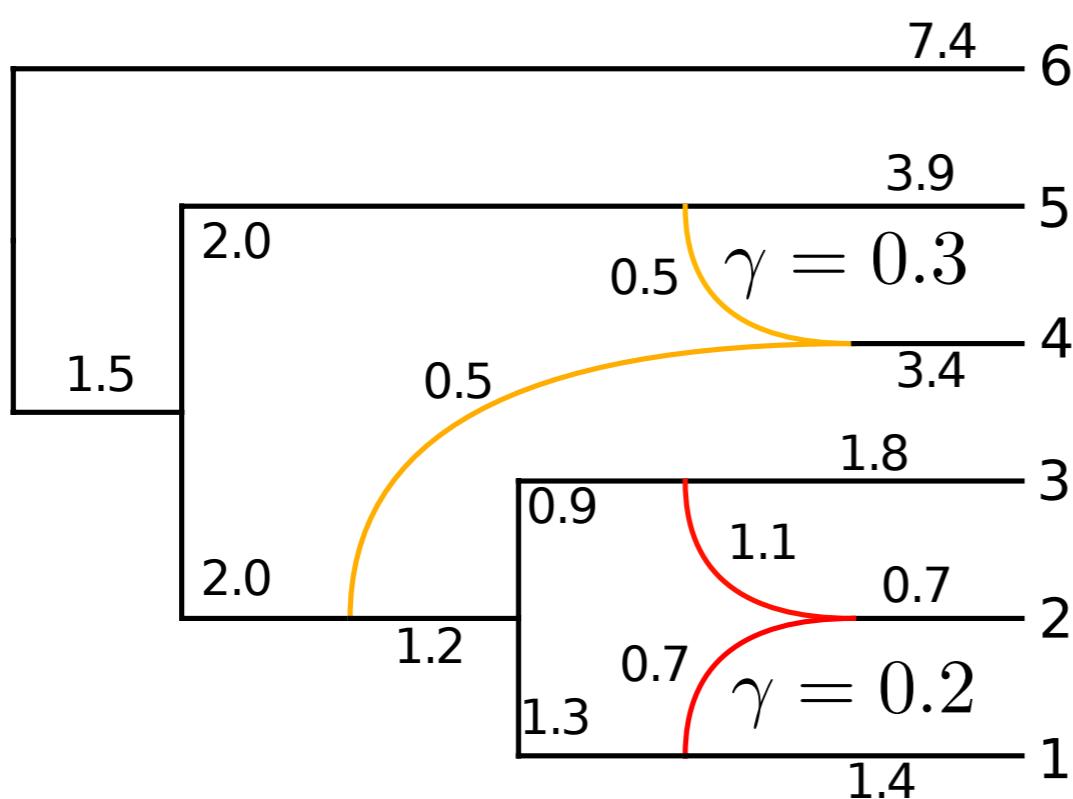
Good diamond  
( $n_0, n_2 \geq 2$ )

Generic Identifiability

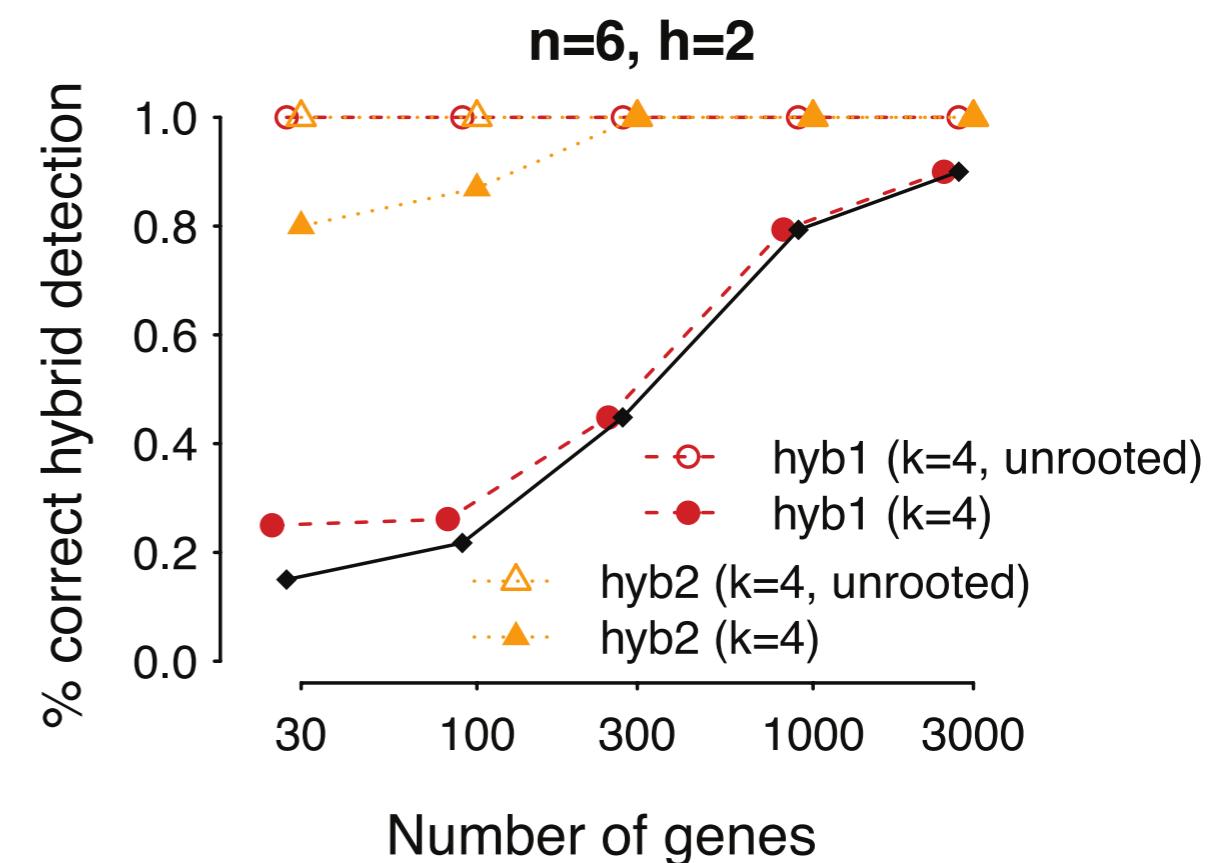
$t_i \in (0, \infty), \gamma \in (0, 1)$

# Identifiability matters: SNaQ performance

Good diamond



Bad diamond



# Challenges

- Network space

- Identifiability

Displayed vs Parental trees  
Level-1 semi-directed networks  
Hybridizations: case by case  
**Missing:** likelihood, level-k semi-directed

- Network comparison

# Challenges

- Network space

K. Huber, V. Moulton, C. Scornavacca,...  
**Missing:** path through tree space, semi-directed

- Identifiability

Displayed vs Parental trees  
Level-1 semi-directed networks  
Hybridizations: case by case  
**Missing:** likelihood, level-k semi-directed

- Network comparison

# Challenges

- Network space

K. Huber, V. Moulton, C. Scornavacca,...  
**Missing:** path through tree space, semi-directed

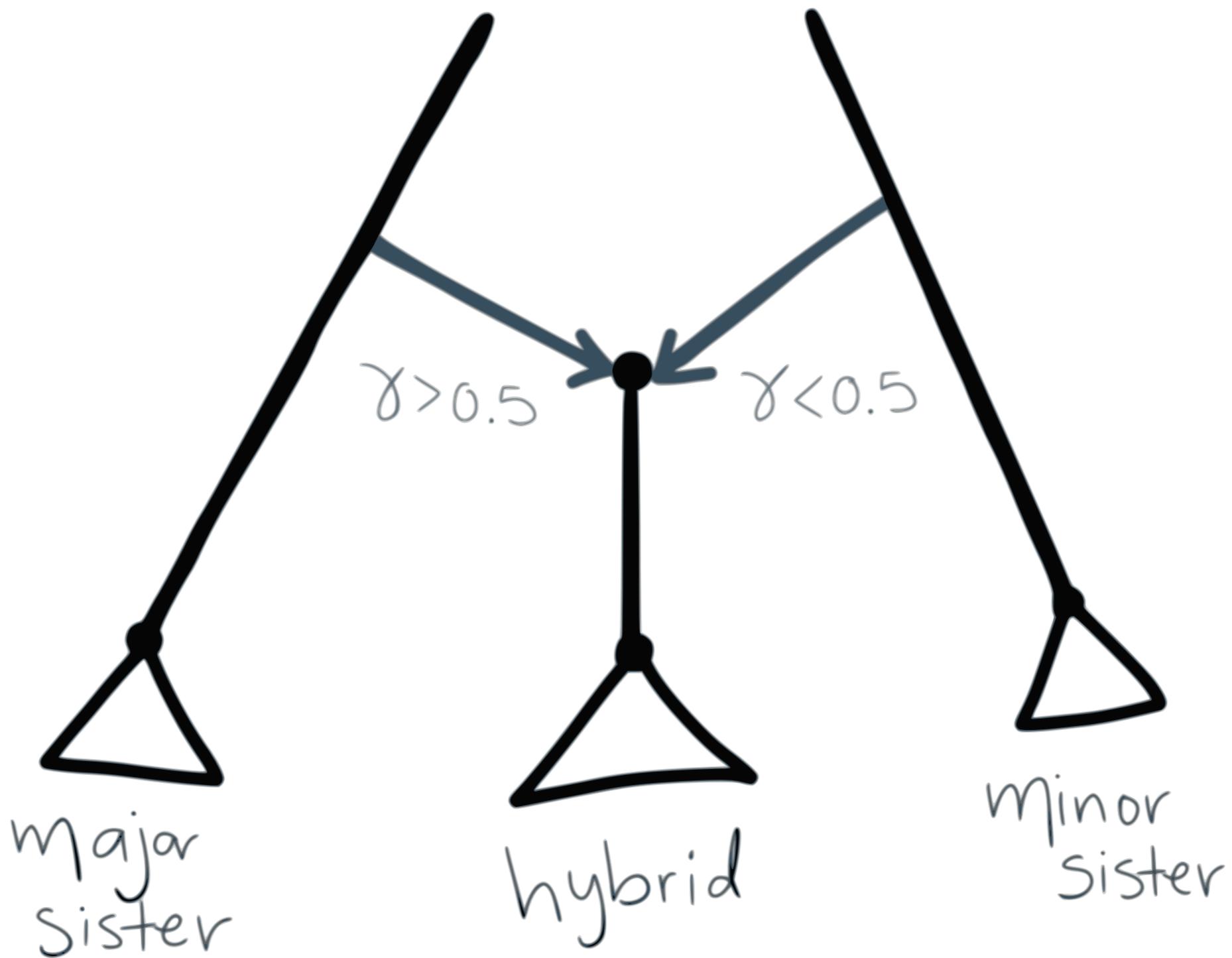
- Identifiability

Displayed vs Parental trees  
Level-1 semi-directed networks  
Hybridizations: case by case  
**Missing:** likelihood, level-k semi-directed

- Network comparison

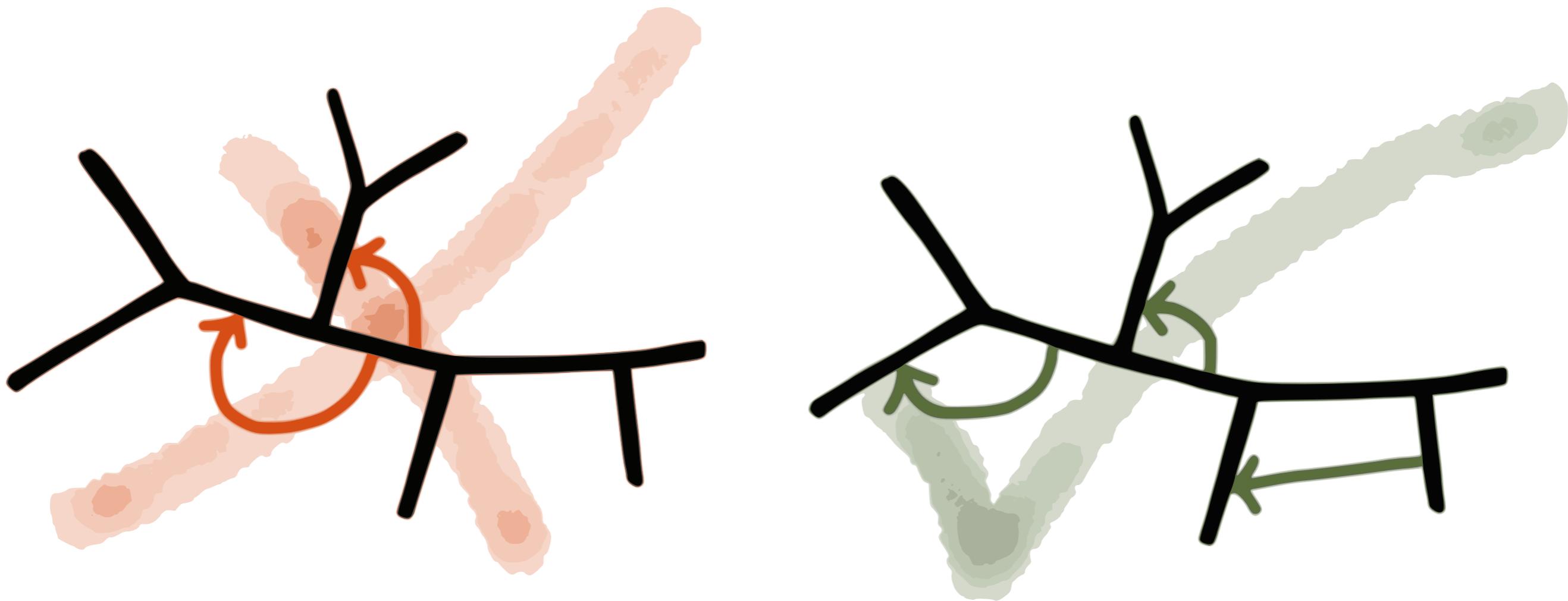
  
**Missing:** distance function  
Hardwired-cluster distance only for rooted networks  
Summary of networks: clades!

# Network summary



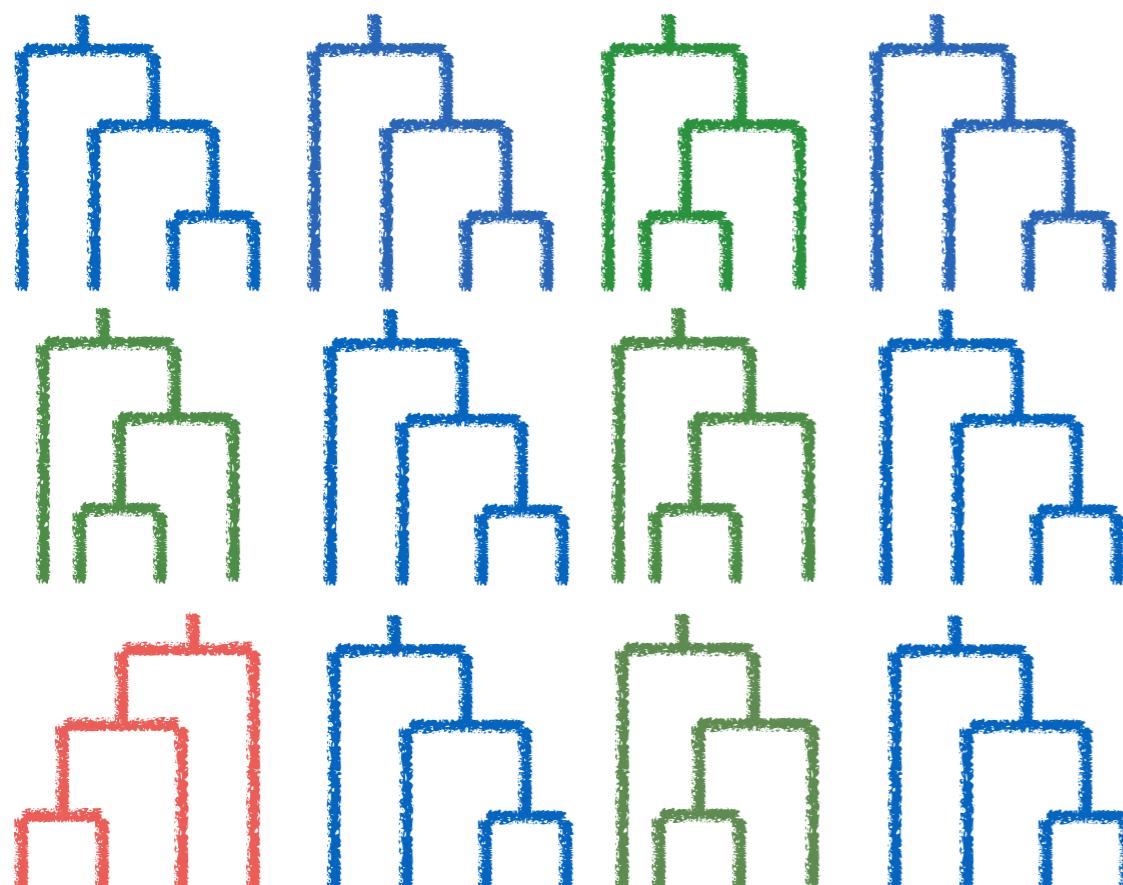
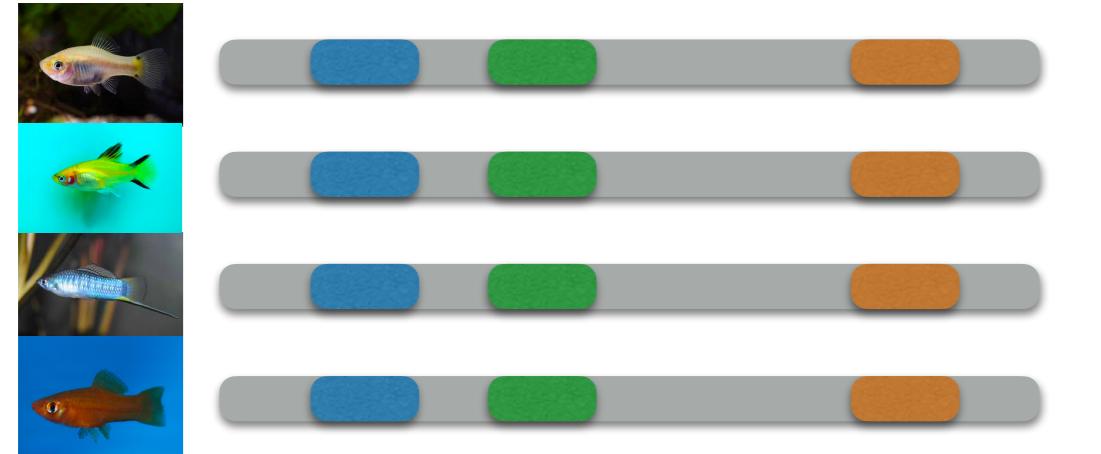
(S.-L. et al, 2017, MBE)

# snaQ limitation: Level-1 networks



# When?

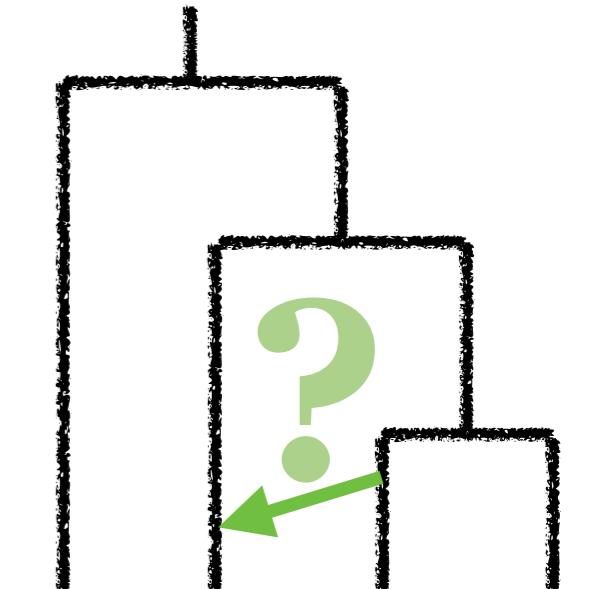
## Phylogenetic network



Data

TICR  
→  
 GitHub

**Goodness-of-fit test**  
Hypothesis test:  
Is a tree a good fit?



<https://github.com/nstenz/TICR>  
(Stenz et al, 2015, Syst Bio)

# PhyloNetworks: analysis for phylogenetic networks

build passing docs stable docs dev codecov 81% coverage 67%

## Overview

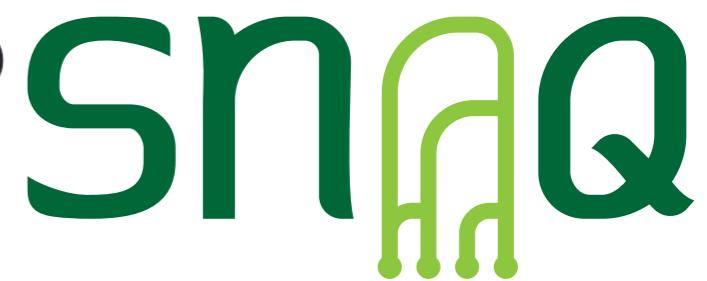


PhyloNetworks is a [Julia](#) package with utilities to:

- read / write phylogenetic trees and networks, in (extended) Newick format. Networks are considered explicit: nodes represent ancestral species. They can be rooted or unrooted.
- manipulate networks: re-root, prune taxa, remove hybrid edges, extract the major tree from a network, extract displayed networks / trees
- compare networks / trees with dissimilarity measures (Robinson-Foulds distance on trees)
- summarize samples of bootstrap networks (or trees) with edge and node support
- estimate species networks from multilocus data (see below)
- phylogenetic comparative methods for continuous trait evolution on species networks / trees



- Step-by-step tutorial
- Online documentation
- Google user group



(S.-L. et al, 2017, MBE)



<https://solislemuslab.github.io/>



@solislemuslab



crsl4

# In-class dynamic

- **Time:** 20 minutes
- **Instructions:** We will go over the PhyloNetworks pipeline which will cover ASTRAL, BUCKy and SNaQ. Create our own reproducible script.
- **Options for you:**
  1. "I think that I can follow the pipeline by myself or with a small group of peers": you should join the Congregate room
  2. "I think I need more one-on-one help to run the commands": you can stay here in the zoom room